

Advanced Algorithms and Models for Computational Biology

-- a machine learning approach

Computational Genomics III: Motif Detection

Eric Xing

Lecture 4, February 1, 2005



Reading: Chap. 1,2, DEKM book

Motifs - Sites - Signals - Domains



- For this lecture, I'll use these terms interchangeably to describe **recurring elements** of interest to us.
- In **PROTEINS** we have: transmembrane domains, coiled-coil domains, EGF-like domains, signal peptides, phosphorylation sites, antigenic determinants, ...
- In **DNA / RNA** we have: enhancers, promoters, terminators, splicing signals, translation initiation sites, centromeres, ...

Example: Gcn4



Regulatory Signals

```

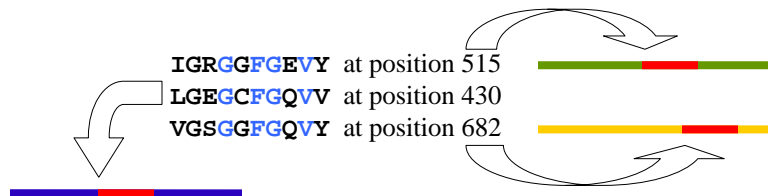
5'- TCTCTCTCCACGGCTAATTAGGTGATCATGAAAAAATGAAAAATTCATGAGAAAAGAGTCAAGACATCGAAACATACAT ...HIS7
5'- ATGGCAGAATCACTTTAAACGTGGCCCCACCCGCTGCACCTGTGCATTTTGTACGTTACTGCGAAATGACTCAACG ...ARO4
5'- CACATCCAACGAATCACCTCACCGTTATCGTGACTCACTTCTTTTCGCATCGCCGAAGTGCCATAAAAAATATTTTTT ...ILV6
5'- TGCGAACAAAAGAGTCAATTACAACGAGGAAATAGAAGAAAAATGAAAAATTTTCGACAAAAATGTATAGTCATTTCTATC ...THR4
5'- ACAAAAGGTACCTTCTGGCCAATCTCACAGATTTAATATAGTAAATGTTCATGCATAAGACTCATCCGGAACATGAAA ...ARO1
5'- ATTGATTGACTCATTTCCTCTGACTACTACCAGTTCAAAATGTTAGAGAAAAATAGAAAAGCAGAAAAAATAAATAA ...HOM2
5'- GGCGCCACAGTCCGCGTTTGGTTATCCGGCTGACTCATTCGACTCTTTTTGGAAAGTGTGGCATGTGCTTCACACA ...PRO3
    
```

Given a collection of genes with common expression,
Can we find the TF-binding site in common?

Motif discovery problem



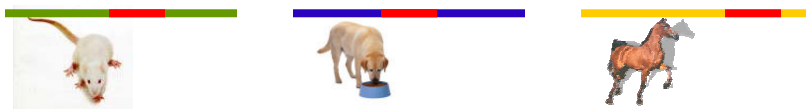
- Given sequences
- Find motif
 - the number of motifs
 - the width of each motif
 - the locations of motif occurrences





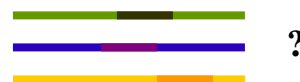
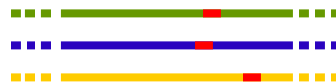
Why find motifs?

- In proteins—may be a critical component
 - Find similarities to known proteins
 - Find important areas of new protein family
- In DNA—may be a *binding site*
 - Discover how the gene expression is regulated



Why is this hard?

- Input sequences are long (thousands or millions of residues)
- Motif may be *subtle*
 - Instances are short.
 - Instances may be only slightly similar.



Characteristics of Regulatory Motifs



```
ATATAAA TT T
CTGATA A CAG
GTGA TCA
AGGGGG ACG
AA AA AA
TTTAA TAA
GAAACGTTGCG
AA TTAATA
TTTAA TAA
CGGACGAG
AAAAATTT
A GA AAAA AA
T TAA T
AA AA AAAA
TTT AA AA
T T T AA AAA
ATAT ATAA
ATTA AAAAT
```

- Tiny
- Highly Variable
- ~Constant Size
 - Because a constant-size transcription factor binds
- Often repeated
- Low-complexity-ish

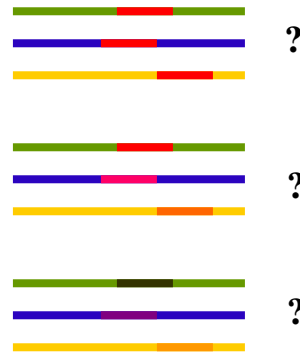
Motif Representation



Measuring similarity



- What counts as a similarity?
- How can such a pattern be searched for?
- Need a *concrete measure* of how good a *motif* is, and how well-matched an *instance* is.



Determinism 1: Consensus Sequences



- σ Factor Promotor **consensus** sequence

	-35	-10
σ^{70}	TTGACA	TATAAT
σ^{28}	CTAAA	CCGATAT

Similarly for σ^{32} , σ^{38} and σ^{54} .

- Consensus sequences have the obvious **limitation**: there is usually some **deviation** from them.

Determinism 2: Regular Expressions



- The characteristic motif of a Cys-Cys-His-His zinc finger DNA binding domain has *regular expression*

C-X(2,4)-C-X(3)-[LIVMFYWC]-X(8)-H-X(3,5)-H

- Here, as in algebra, **X** is unknown. The 29 a.a. sequence of an example domain 1SP1 is as follows, clearly fitting the model.

1SP1:

KKFACPECPKRFMRSDHLSKHIKTHQNKK

Regular Expressions Can Be Limiting



- The regular expression syntax is still **too rigid** to represent many **highly divergent** protein motifs.
- Also, **short** patterns are sometimes insufficient with today's large databases. Even requiring perfect matches you might find many **false positives**. On the other hand some real sites might not be perfect matches.
- We need to go beyond apparently equally likely alternatives, and ranges for gaps. We deal with the former first, having a **distribution at each position**.

Weight Matrix Model (WMM)



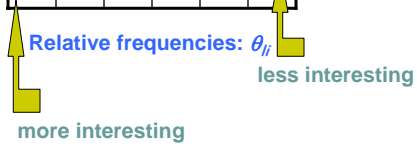
- Weight matrix model (WMM) = Stochastic consensus sequence

A	9	214	63	142	118	8
C	22	7	26	31	52	13
G	18	2	29	38	29	5
T	193	19	124	31	43	216

Counts from 242 known σ^{70} sites

A	.04	.88	.26	.59	.49	.03
C	.09	.03	.11	.13	.21	.05
G	.07	.01	.12	.16	.12	.02
T	.80	.08	.51	.13	.18	.8

Relative frequencies: θ_{li}



- Weight matrices are also known as
 - Position-specific scoring matrices
 - Position-specific probability matrices
 - Position-specific weight matrices
- A motif is *interesting* if it is very different from the background distribution

Weight Matrix Model (WMM)



Weight matrix model (WMM) = Stochastic consensus sequence

A	9	214	63	142	118	8
C	22	7	26	31	52	13
G	18	2	29	38	29	5
T	193	19	124	31	43	216

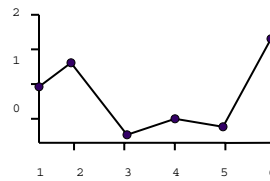
Counts from 242 known σ^{70} sites

A	.04	.88	.26	.59	.49	.03
C	.09	.03	.11	.13	.21	.05
G	.07	.01	.12	.16	.12	.02
T	.80	.08	.51	.13	.18	.89

Relative frequencies: f_{bi}

A	-38	19	1	12	10	-48
C	-15	-38	-8	-10	-3	-32
G	-13	-48	-6	-7	-10	-40
T	17	-32	8	-9	-6	19

$10 \log_2 \theta_{li} / \theta_{oi}$



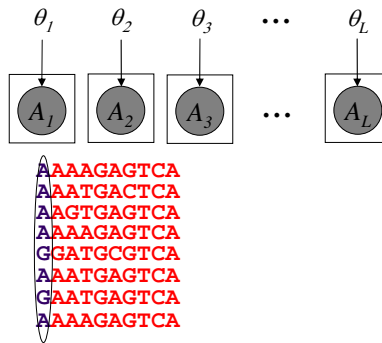
Informativeness: $2 - \sum_i \theta_{li} \log_2 \theta_{li} / \theta_{oi}$

The Product Multinomial (PM) Model

[Lawrence *et al.* Science 1993]



- Positional specific multinomial distribution: $\theta_l = [\theta_{lA}, \dots, \theta_{lC}]^T$



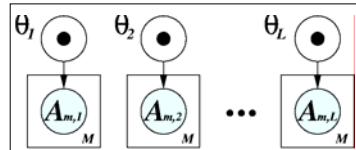
- Position weight matrix (PWM): θ
 - The nucleotide distributions at different positions are independent

More on PM Model



- The PM parameter, $\theta_l = [\theta_{lA}, \dots, \theta_{lC}]^T$, corresponds exactly to the PWM of a motif

	1	2	3	4	5	6	7	8	9	10
A	.750	.875	.875	.375	0	1	0	0	0	1
C	0	0	0	0	0	0	0	.125	1	0
G	.250	.125	.125	0	1	0	.875	0	0	0
T	0	0	0	.635	0	0	0	1	0	0



The nucleotide distributions at different sites are independent !

- The **score** (likelihood-ratio) of a candidate substring: **AAAAGAGTCA**

$$R = \frac{p(x = \{\text{AAAAGAGTCA}\} | \text{PWM})}{p(x = \{\text{AAAAGAGTCA}\} | \text{bk})} = \prod_{l=1}^{10} \frac{p(y_l | \text{PWM})}{p(y_l | \text{bk})} = \prod_{l=1}^{10} \frac{\theta_{l,y_l}}{\theta_{0,y_l}}$$

- Log Likelihood-Ratio: $\text{LR} = \sum_{\text{position } l} \left(\sum_{\text{letter } i} \log_2 \frac{\theta_{l,y_l}}{\theta_{0,y_l}} \right)$

Computational problems for *in silico* motif detection



- Extract a motif model based on (experimentally) identified motifs

```

1: AAAAGATCA
2: AAAAGATCA
. AAGTAGTCA
. AAAAGATCA
. GAGTAGTCA
. AAAAGATCA
. GAGTAGTCA
M: AAAAGATCA
    
```

⇒

AAAAGATCA

- Search for motif instances based on given motif model(s)

```

1: AAAAGATCA
2: AAAAGATCA
. AAGTAGTCA
. AAAAGATCA
. GAGTAGTCA
. AAAAGATCA
. GAGTAGTCA
M: AAAAGATCA
    
```

+

AAAAGATCA

⇒

```

1: AAAAGATCA
2: AAAAGATCA
. AAGTAGTCA
. AAAAGATCA
. GAGTAGTCA
. AAAAGATCA
. GAGTAGTCA
M: AAAAGATCA
    
```

- Uncover novel motifs computationally from genomic sequences

```

1: AAAAGATCA
2: AAAAGATCA
. AAGTAGTCA
. AAAAGATCA
. GAGTAGTCA
. AAAAGATCA
. GAGTAGTCA
M: AAAAGATCA
    
```

⇒

```

1: AAAAGATCA
2: AAAAGATCA
. AAGTAGTCA
. AAAAGATCA
. GAGTAGTCA
. AAAAGATCA
. GAGTAGTCA
M: AAAAGATCA
    
```

+

AAAAGATCA

de novo motif detection

Computational problems for *in silico* motif detection



- Extract a motif model based on (experimentally) identified motifs

Supervised learning

- Search for motif instances based on given motif model(s)

Prediction

- Uncover novel motifs computationally from genomic sequences

Unsupervised learning

Problem definition



Given a collection of promoter sequences s_1, \dots, s_N of genes with common expression

Combinatorial

Motif M: substring $m_1 \dots m_W$

Some of the m_i 's blank

- Find M that occurs in all s_i with $\leq k$ differences
- Or, Find M with smallest total hamming dist

Probabilistic

Motif M: $\{\theta_{ij}\}; 1 \leq i \leq W$

$j = A, C, G, T$

$\theta_{ij} = \text{Prob}[\text{letter } j, \text{ pos } i]$

Find best M, and positions p_1, \dots, p_N in sequences

Use of the matrix to find sites



- Hypothesis:
 - S=site (and independence)
 - R=random (equiprobable, independence)
- Move the matrix along the sequence and score each **window**.
- **Peaks should occur at the true sites.**
- Of course in general any threshold will have some **false positive** and **false negative** rate.

	C	T	A	T	A	A	T	C	
A	-38	19	1	12	10	48			-93
C	-15	-38	-8	-10	-3	-32			
G	-13	-48	-6	-7	-10	-40			
T	17	-32	8	-9	-6	19			

	C	T	A	T	A	A	T	C	
A	-38	19	1	12	10	48			+85
C	-15	-38	-8	-10	-3	-32			
G	-13	-48	-6	-7	-10	-40			
T	17	-32	8	-9	-6	19			

	C	T	A	T	A	A	T	C	
A	-38	19	1	2	10	48			-95
C	-15	-38	-8	-10	-3	-32			
G	-13	-48	-6	-7	-10	-40			
T	17	-32	8	-9	-6	19			

Supervised motif search



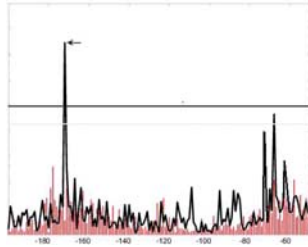
- **Supervised learning**

- Given biologically identified aligned motifs **A**, maximal likelihood estimation:

$$\Theta_{ML} = \arg \max_{\Theta} p(\mathbf{A} | \Theta)$$

- **Application:**

- search for **known** motifs *in silico* from genomic sequences



- Need more more sophisticated search model: HMM?

de novo motif detection



- **Unsupervised learning**

- Given no training examples, predict locations of all instances of **novel** motifs in given sequences, and learn motif models simultaneously.

```

5' - TCTCTCTCCACGGCTAATTAGGTGATCATGAAAAAATGAAAAATTCATGAGAAAAAGATCAGACATCGAAACATACAT ...HIS7
5' - ATGGCAGAATCACTTTAAAAACGTGGCCCCF          FGTACGTTACTGCGAAATGACTCAACG ...ARO4
5' - CACATCCAACGAATCACCTCACCGTTATCCAAATGAGTCA ...CCGAAAGTGCATAAAAAATATTTTTT ...ILV6
5' - TCGGAACAAAAGAGTCAATTACAACGAGGAAATAGAAGAAAAATGAAAAATTTTCGACAAAATGTATAGTCATTTCTATC ...THRA
5' - ACAAAGGTACCTTCCCTGGCCAATCTCACAGATTTAATATAGTAAATTTGTATGCATATGACTCATCCGAAACATGAAA ...ARO1
5' - ATTGATTGACTCATTTCCTCTGACTACTACCAGTTCAAAAATGTTAGAGAAAAATAGAAAAGCAGAAAAATAAATAA ...HOM2
5' - GCGCCACAGTCCGCGTTTGGTTATCCGGCTGACTCATTCTGACTCTTTTTTGGAAAAGTGTGGCATGTGCTTCACACA ...PRO3
    
```

- **Learning algorithms:**

- Expectation Maximization: e.g., MEME
- Gibbs Sampling: e.g., AlignACE, BioProspector
- Advanced models: Bayesian network, Bayesian Markovian models

de novo motif detection



- **Problem setting:**

- Given UTR sequences: $y = \{y_1, \dots, y_N\}$
- Goal: the background model: $\theta_0 = \{\theta_{0,A}, \theta_{0,T}, \theta_{0,G}, \theta_{0,C}\}^t$
and K motif models $\theta_1, \dots, \theta_K$ from y ,
where $\theta^k = \{\theta_{i,j}^k : i = 1, \dots, L_k, j \in \{A, C, G, T\}\}$

- **A missing value problem:**

- The locations of instances of motifs are unknown, thus the aligned motif sequences A_1, \dots, A_K and the background sequence are not available.

Expectation-maximization



For each subsequence of width W
convert subsequence to a matrix

EM [do {
 re-estimate motif occurrences from matrix
 re-estimate matrix model from motif occurrences
} until (matrix model stops changing)
end
select matrix with highest score

Sample DNA sequences



```
>celcg
TAATGTTTGTGCTGGTTTTTGTGGCATCGGGCGAGAATA
GCGCGTGGTGTGAAAGACTGTTTTTTTGGATCGTTTTTCAC
AAAAATGGAAGTCCACAGTCTTGACAG

>ara
GACAAAAACGCGTAACAAAAGTGTCTATAATCACGGCAG
AAAAGTCCACATTGATTATTTGCACGGCGTCACACTTTG
CTATGCCATAGCATTTTTATCCATAAG

>bglr1
ACAAATCCAATAACTTAATTATTGGGATTTGTTATATA
TAACTTTATAAATTCCTAAAATTACACAAAGTTAATAAC
TGTGAGCATGGTCATATTTTTATCAAT

>crp
CACAAAGCGAAAGCTATGCTAAAACAGTCAGGATGCTAC
AGTAATACATTGATGTACTGCATGTATGCAAAGGACGTC
ACATTACCGTGCAGTACAGTTGATAGC
```

Motif occurrences



```
>celcg
taatgtttgtgctggTTTTTgtggcatcgggcgagaata
gcgcgTGGTGTGAAAGACTGTTTTTTGATCGTTTTCAC
aaaaatggaagtccacagtcttgacag

>ara
gacaaaaacgcgtaacaaaagtgtctataatcacggcag
aaaagtccacattgattaTTGCACGGCGTCACactttg
ctatgccatagcatttttatccataag

>bglr1
acaaatccaataacttaattattgggatttgttatata
taactttataaattcctaaaattacacaaagttaataac
TGTGAGCATGGTCATatTTTTatcaat

>crp
caciaagcgaaagctatgctaaaacagtcaggatgctac
agtaatacattgatgtactgcataTGCAAAGGACGTC
ACattaccgtgcagTACAGTTGATAGC
```


Starting point



...gactgTTTT**TTTGATCGTTTTCAC**aaaaatgg...

	T	T	T	G	A	T	C	G	T	T
A	0.17	0.17	0.17	0.17	0.50	...				
C	0.17	0.17	0.17	0.17	0.17					
G	0.17	0.17	0.17	0.50	0.17					
T	0.50	0.50	0.50	0.17	0.17					

This a special initialization scheme, many others scheme, including random starts, are also valid

Re-estimating motif occurrences



TAATGTTTGTGCTGGTTTTTTGTGGCATCGGGCGAGAATA

	T	T	T	G	A	T	C	G	T	T
A	0.17	0.17	0.17	0.17	0.50	...				
C	0.17	0.17	0.17	0.17	0.17					
G	0.17	0.17	0.17	0.50	0.17					
T	0.50	0.50	0.50	0.17	0.17					

Score = 0.50 + 0.17 + 0.17 + 0.17 + 0.17 + ...



Scoring each subsequence

- Score from each sequence the subsequence with maximal score.

Sequence: TGTGCTGGTTTTTGTGGCATCGGGCGAGAATA

Subsequences	Score
TGTGCTGGTTTTTGT	2.95
GTGCTGGTTTTTGTG	4.62
TGCTGGTTTTTGTGG	2.31
GCTGGTTTTTGTGGC	...



Re-estimating motif matrix

- From each sequence, take the substring that has the maximal score
- Align all of them and count:

Occurrences	Counts
TTTGATCGTTTTCAC	A 000132011000040
TTTGCACGGCGTCAC	C 001010300200403
TGTGAGCATGGTCAT	G 020301131130000
TGCAAAGGACGTCAC	T 423001002114001

Adding pseudocounts



Counts		Counts + Pseudocounts
A 000132011000040	→	A 111243122111151
C 001010300200403		C 112121411311514
G 020301131130000		G 131412242241111
T 423001002114001		T 534112113225112

Converting to frequencies



Counts + Pseudocounts										
A 111243122111151										
C 112121411311514										
G 131412242241111										
T 534112113225112										
	T	T	T	G	A	T	C	G	T	T
A	0.13	0.13	0.13	0.25	0.50	...				
C	0.13	0.13	0.25	0.13	0.25					
G	0.13	0.38	0.13	0.50	0.13					
T	0.63	0.38	0.50	0.13	0.13					

Expectation-maximization



EM

```
For each subsequence of width W
  convert subsequence to a matrix
do {
  re-estimate motif occurrences from matrix
  re-estimate matrix model from motif occurrences
} until (matrix model stops changing)
end
select matrix with highest score
```

- **Problem:**
 - This procedure doesn't allow the motifs to move around very much. Taking the max is too brittle.
- **Solution:**
 - Associate with each start site a probability of motif occurrence.

Converting to probabilities



Sequence: TGTGCTGGTTTTTGTGGCATCGGGCGAGAATA

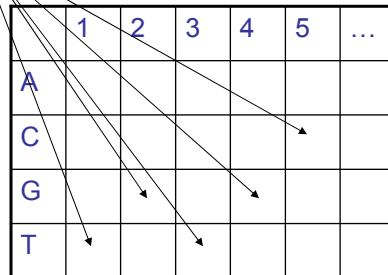
Occurrences	Score	Prob
TGTGCTGGTTTTTGT	2.95	0.023
GTGCTGGTTTTTGTG	4.62	0.037
TGCTGGTTTTTGTGG	2.31	0.018
GCTGGTTTTTGTGGC
Total	128.2	1.000

Computing weighted counts



Occurrences	Prob
TGTGCTGGTTTTTGT	0.023
GTGCTGGTTTTTGTG	0.037
TGCTGGTTTTTGTGG	0.018
GCTGGTTTTTGTGGC	...

Include counts from all subsequences, weighted by the degree to which they match the motif model.



Q. and A.



- **Problem:** How do we estimate counts accurately when we have only a few examples?
 - **Solution:** Use Dirichlet mixture priors.
- **Problem:** Too many possible starting points.
 - **Solution:** Save time by running only one iteration of EM.
- **Problem:** Too many possible widths.
 - **Solution:** Consider widths that vary by $\sqrt{2}$ and adjust motifs afterwards.
- **Problem:** Algorithm assumes exactly one motif occurrence per sequence.
 - **Solution:** Normalize motif occurrence probabilities across all sequences, using a user-specified parameter.



Q. and A.

- **Problem:** The EM algorithm finds only one motif.
 - **Solution:** Probabilistically erase the motif from the data set, and repeat.
- **Problem:** The motif model is too simplistic.
 - **Solution:** Use a two-component mixture model that captures the background distribution. Allow the background model to be more complex.
- **Problem:** The EM algorithm does not tell you how many motifs there are.
 - **Solution:** Compute statistical significance of motifs and stop when they are no longer significant.



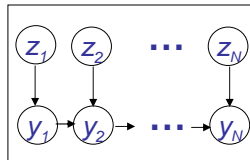
MEME algorithm

```
do
  for (width = min; width *=  $\sqrt{2}$ ; width < max)
    foreach possible starting point
      run 1 iteration of EM
      select candidate starting points
      foreach candidate
        run EM to convergence
      select best motif
      erase motif occurrences
until (E-value of found motif > threshold)
```

What is underlying the EM algorithm? – the statistical foundation



- **A binary indicator model**



$$Z \in \{0, 1\}^N$$

Let:

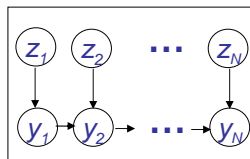
$Y_{n...n+L-1} = \{Y_n, Y_{n+1}, \dots, Y_{n+L-1}\}$: an L-long word starting at position n

$$p(z_n = 1) = \varepsilon, \quad p(z_n = 0) = 1 - \varepsilon$$

$$p(Y_{n...n+L-1} | z_n = 0) = \theta_{0,y_n} \theta_{0,y_{n+1}} \dots \theta_{0,y_{n+L-1}} = \prod_{l=0}^{L-1} \theta_{0,y_{n+l}} = \prod_{l=1}^L \prod_{j=1}^4 \theta_{0,j}^{\delta(Y_{n+l-1}, j)} \quad (\text{background})$$

$$p(Y_{n...n+L-1} | z_n = 1) = \theta_{1,y_n} \theta_{1,y_{n+1}} \dots \theta_{1,y_{n+L-1}} = \prod_{l=1}^L \theta_{l,y_{n+l}} = \prod_{l=1}^L \prod_{j=1}^4 \theta_{l,j}^{\delta(Y_{n+l-1}, j)} \quad (\text{motif seq.})$$

A binary indicator model



$$Z \in \{0, 1\}^N$$

- Complete log-likelihood :

- suppose all words are concatenated into one big sequence of $Y = Y_1 Y_2 \dots Y_N$, with appropriate constraints preventing overlapping and boundaries limits

$$p(Y_n, z_n) = p(y_n | z_n) p(z_n) = p(y_i | \theta_0)^{1-z_n} p(y_n | \theta)^{z_n} \times (1 - \varepsilon)^{1-z_n} \varepsilon^{z_n}$$

$$l_c(\Theta) = \sum_{n=1}^N z_n \left(\sum_{l=1}^L \sum_{j=1}^4 \delta(Y_{n+l-1}, j) \log \theta_{l,j} \right) + \sum_{n=1}^N (1 - z_n) \left(\sum_{j=1}^4 \delta(Y_{n+l-1}, j) \log \theta_{0,j} \right) + |z| \log \varepsilon + (N - |z|) \log(1 - \varepsilon)$$



The Maximal likelihood approach

- Maximize expected likelihood, in iteration of two steps:

Expectation:

Find expected value of complete log likelihood:

$$E[\log \mathcal{P}(Y_1 \dots Y_n, Z | \theta, \theta_0, \varepsilon)]$$

Maximization:

Maximize the expected complete likelihood over $\theta, \theta_0, \varepsilon$



Expectation Maximization: E-step

- Expectation:

Find expected value of log likelihood:

$$\begin{aligned} \langle l_c(\Theta) \rangle = & \sum_{n=1}^N \langle z_n \rangle \left(\sum_{l=1}^L \sum_{j=1}^4 \delta(y_{n+l-1}, j) \log \theta_{l,j} \right) + \sum_{n=1}^N (1 - \langle z_n \rangle) \left(\sum_{j=1}^4 \delta(y_{n+l-1}, j) \log \theta_{0,j} \right) \\ & + \sum_{n=1}^N \langle z_n \rangle \log \varepsilon + \left(N - \sum_{n=1}^N \langle z_n \rangle \right) \log(1 - \varepsilon) \end{aligned}$$

- where the expected value of Z can be computed as follows:

$$\langle z_i \rangle = p(z_i = 1 | Y) = \frac{\varepsilon p(y_i | \theta)}{\varepsilon p(y_i | \theta) + (1 - \varepsilon) p(y_i | \theta_0)}$$

- recall the weights for each substring in the MEME algorithm

Expectation Maximization: M-step



- Maximization:

Maximize expected value over θ and ε independently

For ε , this is easy:

$$\varepsilon^{NEW} = \arg \max_{\varepsilon} \sum_{n=1}^N \langle z_n \rangle \log \varepsilon + (N - \sum_{n=1}^N \langle z_n \rangle) \log(1 - \varepsilon) = \frac{\sum_{n=1}^N \langle z_n \rangle}{N}$$

Expectation Maximization: M-step



- For $\theta = (\theta, \theta_0)$, define

$c_{i,j} = E[\text{\# times letter } j \text{ appears in motif position } i]$

$c_{0,j} = E[\text{\# times letter } j \text{ appears in background}]$

- $c_{i,j}$ values are calculated easily from $E[Z]$ values

It easily follows:

$$\theta_{i,j}^{NEW} = \frac{c_{i,j}}{\sum_{j=1}^A c_{i,j}} \quad \theta_{0,j}^{NEW} = \frac{c_{0,j}}{\sum_{j=1}^A c_{0,j}}$$

to not allow any 0's, add pseudocounts



Initial Parameters Matter!

- Consider the following “artificial” example:

x^1, \dots, x^N contain:

- 2^{12} patterns on {A, T}: A...AA, A...AT, ..., T...TT
- 2^{12} patterns on {C, G}: C...CC, C...CG, ..., G...GG
- $D \ll 2^{12}$ occurrences of 12-mer ACTGACTGACTG

- Some local maxima:

$$\varepsilon \approx \frac{1}{2}; \quad B = \frac{1}{2}C, \frac{1}{2}G; \quad M_i = \frac{1}{2}A, \frac{1}{2}T, \quad i = 1, \dots, 12$$

- very bad !

$$\varepsilon \approx D/2^{L+1}; \quad B = \frac{1}{4}A, \frac{1}{4}C, \frac{1}{4}G, \frac{1}{4}T; \\ M_1 = 100\% A, M_2 = 100\% C, M_3 = 100\% T, \text{ etc.}$$

- the correct solution !



Overview of EM Algorithm

1. Initialize parameters $\theta = (\theta, \theta_0), :$
 - Try different values of ε , say, from $N^{-1/2}$ up to $1/(2L)$
2. Repeat:
 - a. Expectation
 - b. Maximization
3. Until change in $\theta = (\theta, \theta_0)$, falls below δ
4. Report results for several “good” ε

Overview of EM Algorithm



- One iteration running time: $O(NL)$
 - Usually need $< N$ iterations for convergence, and $< N$ starting points.
 - Overall complexity: unclear
- EM is a local optimization method
- Initial parameters matter
- MEME: Bailey and Elkan, ISMB 1994.