

# Learning Sequence Motif Models Using Expectation Maximization (EM)

BMI/CS 776

[www.biostat.wisc.edu/bmi776/](http://www.biostat.wisc.edu/bmi776/)

Spring 2011

Mark Craven

[craven@biostat.wisc.edu](mailto:craven@biostat.wisc.edu)

## Goals for Lecture

the key concepts to understand are the following

- the motif finding problem
- using EM to address the motif-finding problem
- the OOPS and ZOOPS models

# Sequence Motifs

- what is a sequence *motif* ?
  - a sequence pattern of biological significance
- examples
  - DNA sequences corresponding to protein binding sites
  - protein sequences corresponding to common functions or conserved pieces of structure

## Sequence Motifs Example

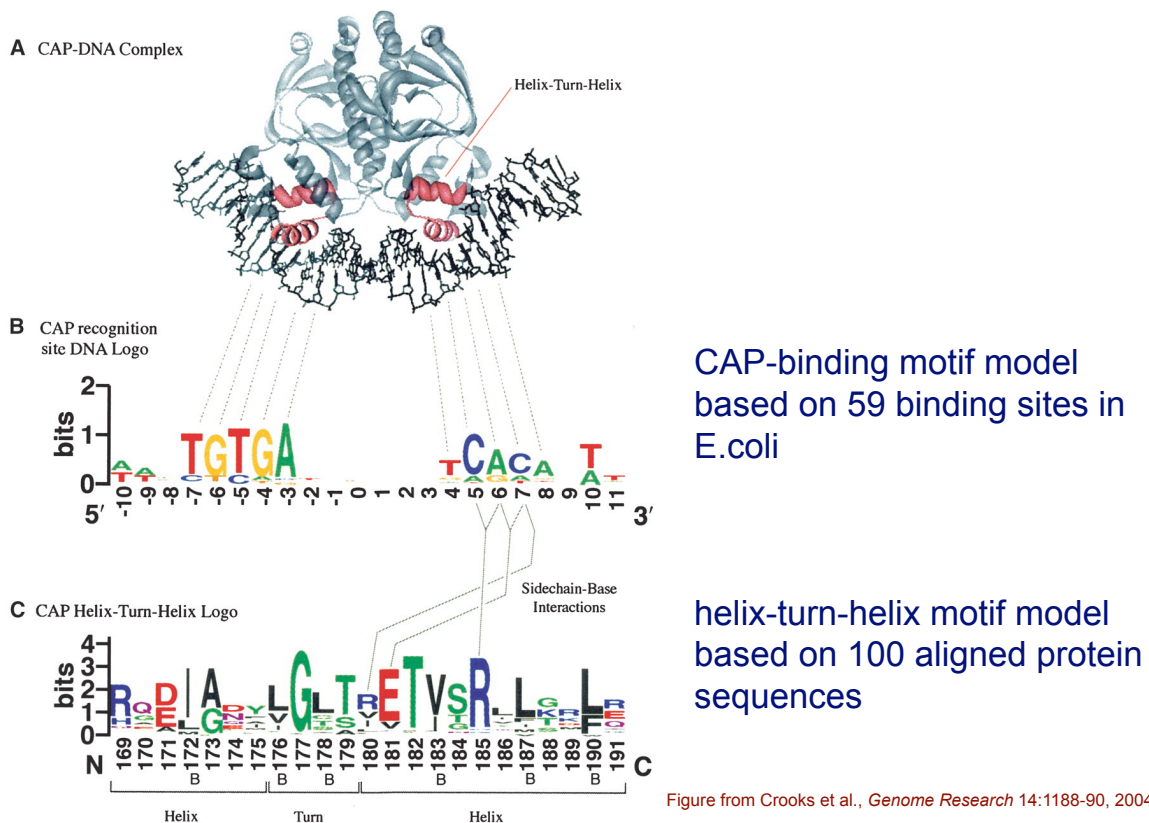


Figure from Crooks et al., *Genome Research* 14:1188-90, 2004.

# The Motif Model Learning Task

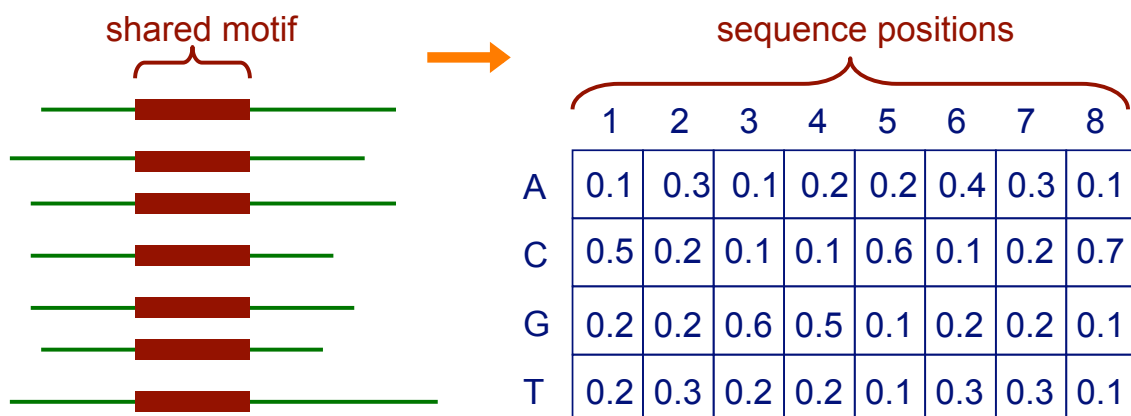
**given:** a set of sequences that are thought to contain an unknown motif of interest

**do:**

- infer a model of the motif
- predict the locations of the motif in the given sequences

## Motifs and *Profile Matrices* (a.k.a. *Position Weight Matrices*)

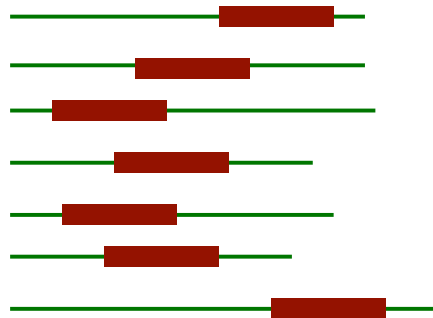
- given a set of aligned sequences, it is straightforward to construct a profile matrix characterizing a motif of interest



- each element represents the probability of given character at a specified position

# Motifs and Profile Matrices

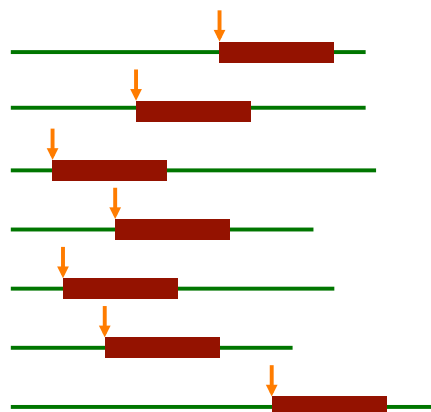
- How can we construct the profile if the sequences aren't aligned?
- In the typical case we don't know what the motif looks like.



## The EM Approach

[Lawrence & Reilly, 1990; Bailey & Elkan, 1993, 1994, 1995]

- EM is a family of algorithms for learning probabilistic models in problems that involve *hidden state*
- in our problem, the hidden state is where the motif starts in each training sequence




# Representing Motifs in MEME

- a motif is
  - assumed to have a fixed width,  $W$
  - represented by a matrix of probabilities:  $p_{c,k}$   
represents the probability of character  $c$  in column  $k$
- also represent the “background” (i.e. sequence outside the motif):  $p_{c,0}$  represents the probability of character  $c$  in the background

# Representing Motifs in MEME

- example: a motif model of length 3

		0	1	2	3
$p =$	A	0.25	0.1	0.5	0.2
	C	0.25	0.4	0.2	0.1
	G	0.25	0.3	0.1	0.6
	T	0.25	0.2	0.2	0.1



background      motif positions

# Representing Motif Starting Positions in MEME

- the element  $Z_{i,j}$  of the matrix  $Z$  represents the probability that the motif starts in position  $j$  in sequence  $i$
- example: given DNA sequences of length 6, where  $W=3$

Diagram showing four DNA sequences of length 6. The first three sequences have a motif of length 3 ( $W=3$ ) starting at position 1, 2, and 3 respectively. The fourth sequence has a motif starting at position 4. The matrix  $Z$  below shows the probability of the motif starting at each position for each sequence.

		1	2	3	4
	seq1	0.1	0.1	0.2	0.6
	seq2	0.4	0.2	0.1	0.3
	seq3	0.3	0.1	0.5	0.1

## Likelihood of a Sequence Given a Motif Starting Position



$$P(X_i | Z_{i,j} = 1, p) = \underbrace{\prod_{k=1}^{j-1} p_{c_k, 0}}_{\text{before motif}} \underbrace{\prod_{k=j}^{j+W-1} p_{c_k, k-j+1}}_{\text{motif}} \underbrace{\prod_{k=j+W}^L p_{c_k, 0}}_{\text{after motif}}$$

$X_i$  is the  $i$  th sequence

$Z_{i,j}$  is 1 if motif starts at position  $j$  in sequence  $i$

$c_k$  is the character at position  $k$  in sequence  $i$

## Likelihood Example

$$X_i = \text{G C } \boxed{\text{T G T}} \text{ A G}$$

$$p = \begin{array}{cc} & \begin{array}{c} 0 \quad 1 \quad 2 \quad 3 \end{array} \\ \begin{array}{c} \text{A} \\ \text{C} \\ \text{G} \\ \text{T} \end{array} & \begin{array}{cccc} 0.25 & 0.1 & 0.5 & 0.2 \\ 0.25 & 0.4 & 0.2 & 0.1 \\ 0.25 & 0.3 & 0.1 & 0.6 \\ 0.25 & 0.2 & 0.2 & 0.1 \end{array} \end{array}$$

$$P(X_i \mid Z_{i3} = 1, p) =$$

$$p_{\text{G},0} \times p_{\text{C},0} \times p_{\text{T},1} \times p_{\text{G},2} \times p_{\text{T},3} \times p_{\text{A},0} \times p_{\text{G},0} = \\ 0.25 \times 0.25 \times 0.2 \times 0.1 \times 0.1 \times 0.25 \times 0.25$$

## Basic EM Approach

given: length parameter  $W$ , training set of sequences

$t=0$

set initial values for  $p^{(0)}$

do

$++t$

  re-estimate  $Z^{(t)}$  from  $p^{(t-1)}$  (E –step)

  re-estimate  $p^{(t)}$  from  $Z^{(t)}$  (M-step)

until change in  $p^{(t)} < \epsilon$

return:  $p^{(t)}, Z^{(t)}$

## The E-step: Estimating Z

- to estimate the starting positions in Z at step  $t$

$$Z_{i,j}^{(t)} = \frac{P(X_i | Z_{i,j} = 1, p^{(t-1)})P(Z_{i,j} = 1)}{\sum_{k=1}^{L-W+1} P(X_i | Z_{i,k} = 1, p^{(t-1)})P(Z_{i,k} = 1)}$$

- this comes from Bayes' rule applied to

$$P(Z_{i,j} = 1 | X_i, p^{(t-1)})$$

## The E-step: Estimating Z

- assume that it is equally likely that the motif will start in any position

$$Z_{i,j}^{(t)} = \frac{P(X_i | Z_{i,j} = 1, p^{(t-1)}) \cancel{P(Z_{i,j} = 1)}}{\sum_{k=1}^{L-W+1} P(X_i | Z_{i,k} = 1, p^{(t-1)}) \cancel{P(Z_{i,k} = 1)}}$$



## Example: Estimating $Z$

$$X_i = \text{G C T G T A G}$$

$$p = \begin{array}{cc} & \begin{matrix} 0 & 1 & 2 & 3 \end{matrix} \\ \begin{matrix} \text{A} \\ \text{C} \\ \text{G} \\ \text{T} \end{matrix} & \begin{bmatrix} 0.25 & 0.1 & 0.5 & 0.2 \\ 0.25 & 0.4 & 0.2 & 0.1 \\ 0.25 & 0.3 & 0.1 & 0.6 \\ 0.25 & 0.2 & 0.2 & 0.1 \end{bmatrix} \end{array}$$

$$Z_{i,1} = 0.3 \times 0.2 \times 0.1 \times 0.25 \times 0.25 \times 0.25 \times 0.25$$

$$Z_{i,2} = 0.25 \times 0.4 \times 0.2 \times 0.6 \times 0.25 \times 0.25 \times 0.25$$

⋮

- then normalize so that  $\sum_{j=1}^{L-W+1} Z_{i,j} = 1$

## The M-step: Estimating $p$

- recall  $p_{c,k}$  represents the probability of character  $c$  in position  $k$ ; values for  $k=0$  represent the background

$$p_{c,k}^{(t)} = \frac{n_{c,k} + d_{c,k}}{\sum_b (n_{b,k} + d_{b,k})}$$

pseudo-counts

$$n_{c,k} = \begin{cases} \sum_i \sum_{\{j | X_{i,j+k-1} = c\}} Z_{i,j}^{(t)} & k > 0 \\ n_c - \sum_{j=1}^W n_{c,j} & k = 0 \end{cases}$$

sum over positions where  $c$  appears

total # of  $c$ 's in data set  $\rightarrow n_c$

## Example: Estimating $p$

**A C A G C A**

$$Z_{1,1} = 0.1, Z_{1,2} = 0.7, Z_{1,3} = 0.1, Z_{1,4} = 0.1$$

**A G G C A G**

$$Z_{2,1} = 0.4, Z_{2,2} = 0.1, Z_{2,3} = 0.1, Z_{2,4} = 0.4$$

**T C A G T C**

$$Z_{3,1} = 0.2, Z_{3,2} = 0.6, Z_{3,3} = 0.1, Z_{3,4} = 0.1$$

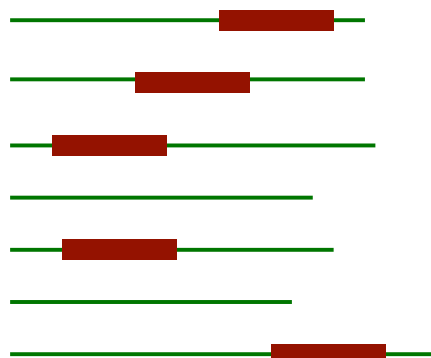
$$p_{A,1} = \frac{Z_{1,1} + Z_{1,3} + Z_{2,1} + Z_{3,3} + 1}{Z_{1,1} + Z_{1,2} + \dots + Z_{3,3} + Z_{3,4} + 4}$$

$$p_{C,2} = \frac{Z_{1,1} + Z_{1,4} + Z_{2,3} + Z_{3,1} + 1}{Z_{1,1} + Z_{1,2} + \dots + Z_{3,3} + Z_{3,4} + 4}$$

⋮

## The ZOOPS Model

- the approach as we've outlined it, assumes that each sequence has exactly one motif occurrence per sequence; this is the OOPS model
- the ZOOPS model assumes zero or one occurrences per sequence



## E-step in the ZOOPS Model

- we need to consider another alternative: the  $i$ th sequence doesn't contain the motif
- we add another parameter (and its relative)

$\lambda$

- prior probability that any position in a sequence is the start of a motif

$\gamma = (L - W + 1)\lambda$

- prior probability of a sequence containing a motif

## E-step in the ZOOPS Model

$$Z_{i,j}^{(t)} = \frac{P(X_i | Z_{i,j} = 1, p^{(t-1)}) \lambda^{(t-1)}}{P(X_i | Q_i = 0, p^{(t-1)})(1 - \gamma^{(t-1)}) + \sum_{k=1}^{L-W+1} P(X_i | Z_{i,k} = 1, p^{(t-1)}) \lambda^{(t-1)}}$$

- $Q_i$  is a random variable for which  $Q_i = 1$  if sequence  $X_i$  contains a motif,  $Q_i = 0$  otherwise

$$P(Q_i = 1) = \sum_{j=1}^{L-W+1} Z_{i,j}^{(t-1)}$$

$$P(X_i | Q_i = 0, p^{(t-1)}) = \prod_{j=1}^L p_{c_j,0}^{(t-1)}$$

## M-step in the ZOOPS Model

- update  $p$  same as before
- update  $\gamma$  as follows:

$$\gamma^{(t)} \equiv (L - W + 1)\lambda^{(t)} = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{L-W+1} Z_{i,j}^{(t)}$$

## Extensions to the Basic EM Approach in MEME

- varying the approach (TCM model) to assume *zero or more* motif occurrences per sequence
- choosing the width of the motif
- finding multiple motifs in a group of sequences
- ✓ choosing good starting points for the parameters
- ✓ using background knowledge to bias the parameters

## Starting Points in MEME

- EM is susceptible to local maxima, so it's a good idea to try multiple starting points
- insight: motif must be similar to *some* subsequence in data set
- for every distinct subsequence of length  $W$  in the training set
  - derive an initial  $p$  matrix from this subsequence
  - run EM for 1 iteration
- choose motif model (i.e.  $p$  matrix) with highest likelihood
- run EM to convergence

## Using Subsequences as Starting Points for EM

- set values matching letters in the subsequence to some value  $\pi$
- set other values to  $(1 - \pi)/(M - 1)$  where  $M$  is the length of the alphabet
- example: for the subsequence TAT with  $\pi = 0.5$

$$p = \begin{array}{ccccc} & & 1 & 2 & 3 \\ \begin{array}{c} \text{A} \\ \text{C} \\ \text{G} \\ \text{T} \end{array} & = & \begin{array}{ccc} 0.17 & 0.5 & 0.17 \\ 0.17 & 0.17 & 0.17 \\ 0.17 & 0.17 & 0.17 \\ 0.5 & 0.17 & 0.5 \end{array} \end{array}$$