

Towards a more adequate motif scoring function

CSRE element in yeast

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
C	G/C	G/T	T/A	C/T	G/C	C/G	A	T	G/T	C/G	A	T	C/T	C/T	G/T

5 strongly conserved positions and eleven weakly conserved positions each of which features two nucleotides with similar frequencies.

NF-kB

COUNT(Motifs)	A:	2	2	0	0	0	0	9	1	1	1	3	0
	C:	1	6	0	0	0	0	0	4	1	2	4	6
	G:	0	0	10	10	9	9	1	0	0	0	0	0
	T:	7	2	0	0	1	1	0	5	8	7	3	4

1	2	3	4	5	6	7	8	9	10	11	12
T	C	G	G	G	G	A	T/C	T	T	C	C/T

last column is *more conserved* than the second column and should receive a lower score

Entropy corresponds to a probability distribution, it is a measure of the uncertainty of a probability distribution.

		↓					↓					↓	
COUNT(<i>Motifs</i>)	A:	2	2	0	0	0	0	9	1	1	1	3	0
	C:	1	6	0	0	0	0	0	4	1	2	4	6
	G:	0	0	10	10	9	9	1	0	0	0	0	0
	T:	7	2	0	0	1	1	0	5	8	7	3	4

→ $0.2 \log_2 0.2 + 0.6 \log_2 0.6 + 0 \log_2 0 + 0.2 \log_2 0.2 = 1.371$

→ $0 \log_2 0 + 0.6 \log_2 0.6 + 0 \log_2 0 + 0.4 \log_2 0.4 = 0.971$

→ $0 \log_2 0 + 0 \log_2 0 + 0.9 \log_2 0.9 + 0.1 \log_2 0.1 = 0.467$

The more conserved the columns the smaller its entropy

The motif finding problem - *Score*

With matrices

Motif Finding Problem

given a collection of strings, find a set of k-mers, one from each string, that minimises the score of the resulting motif.

input: A collection of strings DNA and a integer k

output: A collection Motif of k-mers, one from each string in DNA, minimising SCORE(Motifs) among all possible choices of k-mers

	T	C	G	G	G	G	g	T	T	T	t	t	
	c	C	G	G	t	G	A	c	T	T	a	C	
	a	C	G	G	G	G	A	T	T	T	t	C	
	T	t	G	G	G	G	A	c	T	T	t	t	
	a	a	G	G	G	G	A	c	T	T	C	C	
	T	t	G	G	G	G	A	c	T	T	C	C	
	T	C	G	G	G	G	A	T	T	c	a	t	
	T	C	G	G	G	G	A	T	T	c	C	t	
	T	a	G	G	G	G	A	a	c	T	a	C	
	T	C	G	G	G	t	A	T	a	a	C	C	
SCORE(Motifs)	3	4	0	0	1	1	1	5	2	3	6	4	= 30

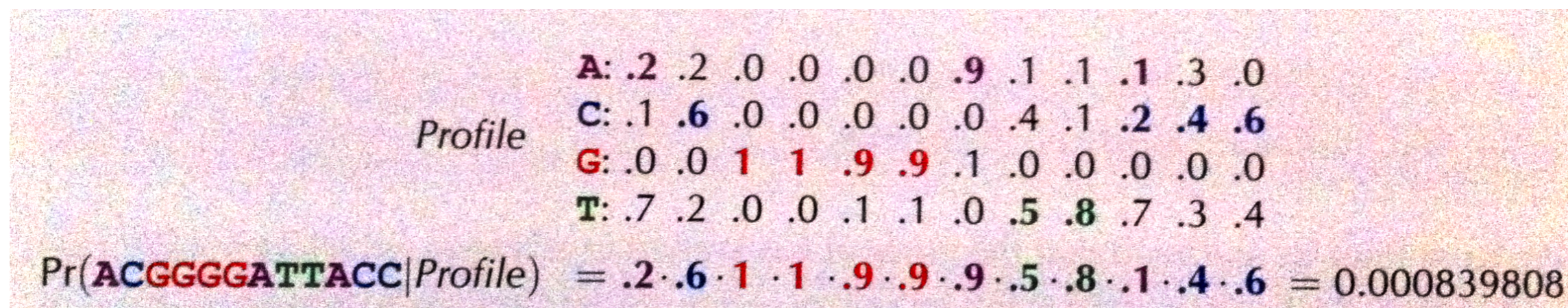
BruteForce algorithm, for every possible kmer from a DNA it is necessary compute the score value and returns those having the minimum score

The motif finding problem - Profile

Motifs be a collection of k-mers taken from t strings of DNA

The *entropy* of each column of a PWM is like a four-side dice:

- A 0.2 A generated with probability 0.2
- C 0.1 C generated with probability 0.1
- G 0.0 G generated with probability 0.0
- T 0.7 T generated with probability 0.7



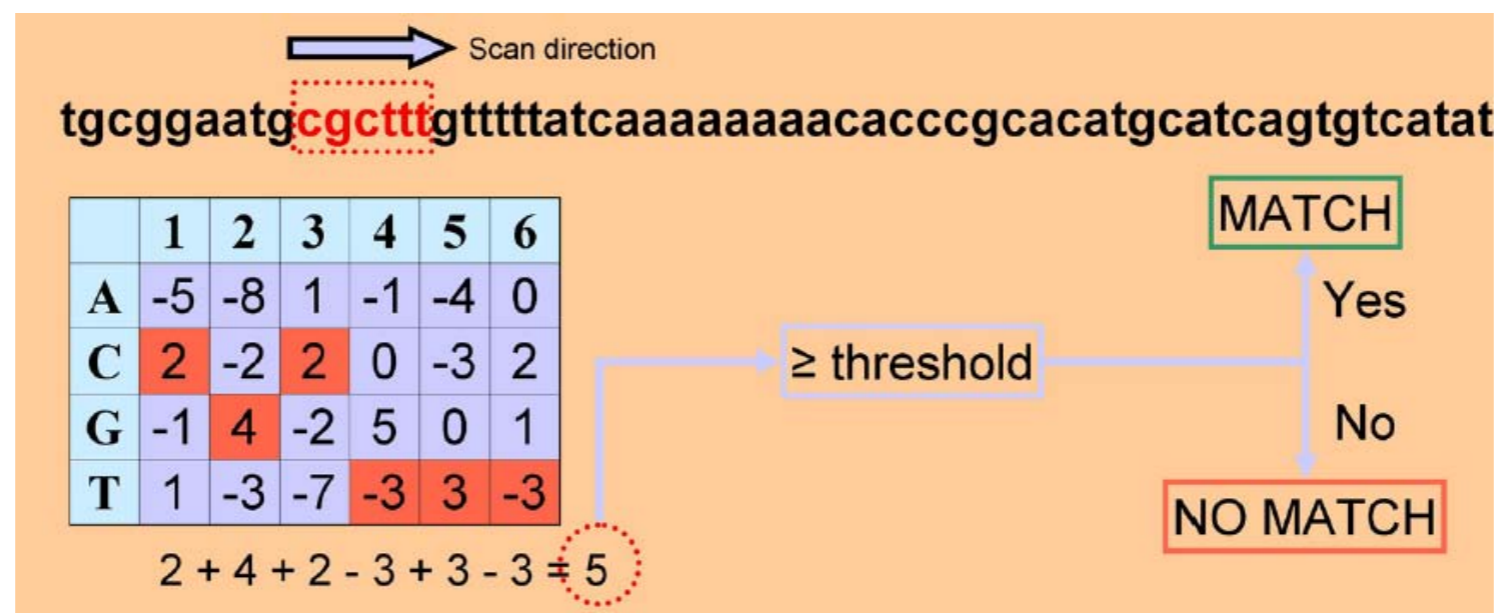
Profile-most Probable k-mer Problem

find a Profile-most probable k-mer in a string

input: A string DNA, a integer k, and a 4 X k matrix Profile

output: A profile-most probable k-mer in DNA

Scoring a Sequence



Courtesy of Kenzie MacIsaac and Ernest Fraenkel. Used with permission. MacIsaac, Kenzie, and Ernest Fraenkel. "Practical Strategies for Discovering Regulatory DNA Sequence Motifs." *PLoS Computational Biology* 2, no. 4 (2006): e36.

Common threshold = 60% of maximum score

MacIsaac & Fraenkel (2006) PLoS Comp Bio

Scoring A Sequence

To score a sequence, we compare to a null model

$$Score = \log \frac{P(S | PFM)}{P(S | B)} = \log \prod_{i=1}^N \frac{P_i(S_i | PFM)}{P(S_i | B)} = \sum_{i=1}^N \underbrace{\log \frac{P_i(S_i | PFM)}{P(S_i | B)}}_{}$$

Equivalent Motif Finding Problem

given a collection of strings, find a pattern and a collection of k-mers (one for each string) that minimises the distance between all possible patterns and all possible collections of k-mers.

input: A collection of strings DNA and a integer k

output: A k-mer pattern and a collection of k-mers Motifs, one from each string in DNA, minimising $d(\text{Pattern}, \text{Motifs})$ among all possible choices of Pattern and Motifs

$d(\text{pattern}, \text{motifs}) = \text{HAMMING DISTANCE}(\text{Pattern}, \text{Motif})$

BUT given a Pattern, we do not need to explore all possible collection of Motifs in order to minimise $d(\text{Pattern}, \text{Motifs})$

$$d(\text{Pattern}, \text{Text}) = \min_{\text{all } k\text{-mers } \text{Pattern}' \text{ in Text}} \text{HAMMINGDISTANCE}(\text{Pattern}, \text{Pattern}').$$

$$d(\text{GATTCTCA}, \text{gcaaaaGACGCTGA} \text{ccaa}) = 3.$$

$$\text{MOTIF}(\text{GATTCTCA}, \text{gcaaaaGACGCTGA} \text{ccaa}) = \text{GACGCTGA}.$$

Median string Problem

find a median string

input: A collection of strings DNA and a integer k

output: A k -mer pattern minimising $d(\text{Pattern}, \text{DNA})$ among all k -mers Pattern

```
MEDIANSTRING(Dna, k)
  distance  $\leftarrow$   $\infty$ 
  for each k-mer Pattern from AA...AA to TT...TT
    if distance  $>$   $d(\text{Pattern}, \text{Dna})$ 
      distance  $\leftarrow$   $d(\text{Pattern}, \text{Dna})$ 
      Median  $\leftarrow$  Pattern
  return Median
```

Greedy motif search

Greedy algorithm select the most attractive alternative at each iteration

```
GREEDYMOTIFSEARCH(Dna, k, t)  
  BestMotifs ← motif matrix formed by first k-mers in each string from Dna  
  for each k-mer Motif in the first string from Dna  
    Motif1 ← Motif  
    for i = 2 to t  
      form Profile from motifs Motif1, ..., Motifi-1  
      Motifi ← Profile-most probable k-mer in the i-th string in Dna  
  Motifs ← (Motif1, ..., Motift)  
  if SCORE(Motifs) < SCORE(BestMotifs)  
    BestMotifs ← Motifs  
return BestMotifs
```



ttACCTtaac
gATGTctgtc
acgGCGTtag
ccctaACGAg
cgtcagAGGT

If the algorithm has chosen ACCT from the first string, what is the profile?

Laplace's rule of succession

	A:	.2	.2	.0	.0	.0	.0	.9	.1	.1	.1	.3	.0		
	C:	.1	.6	.0	.0	.0	.0	.0	.4	.1	.2	.4	.6		
	G:	.0	.0	1	1	.9	.9	.1	.0	.0	.0	.0	.0		
	T:	.7	.2	.0	.0	.1	.1	.0	.5	.8	.7	.3	.4		
<i>Profile</i>															
Pr(TCGTGGATTCC <i>Profile</i>)		=	.7	.6	1	.0	.9	.9	.9	.5	.8	.7	.4	.6	= 0

To improve the unfair scoring, often substitute zeros with small number called **pseudo counts**

	T	A	A	C
<i>Motifs</i>	G	T	C	T
	A	C	T	A
	A	G	G	T

	A:	2	1	1	1
<i>COUNT(Motifs)</i>	C:	0	1	1	1
	G:	1	1	1	0
	T:	1	1	1	2

	2/4	1/4	1/4	1/4
<i>PROFILE(Motifs)</i>	0	1/4	1/4	1/4
	1/4	1/4	1/4	0
	1/4	1/4	1/4	2/4

Laplace's rule of succession adds 1 to each element to COUNT(motif)

	A:	2+1	1+1	1+1	1+1
<i>COUNT(Motifs)</i>	C:	0+1	1+1	1+1	1+1
	G:	1+1	1+1	1+1	0+1
	T:	1+1	1+1	1+1	2+1

	3/8	2/8	2/8	2/8
<i>PROFILE(Motifs)</i>	1/8	2/8	2/8	2/8
	2/8	2/8	2/8	1/8
	2/8	2/8	2/8	3/8

Greedy motif search

Greedy algorithm select the most attractive alternative at each iteration

```
GREEDYMOTIFSEARCH(Dna, k, t)
  BestMotifs ← motif matrix formed by first k-mers in each string from Dna
  for each k-mer Motif in the first string from Dna
    Motif1 ← Motif
    for i = 2 to t
      form Profile from motifs Motif1, ..., Motifi-1}
      Motifi ← Profile-most probable k-mer in the i-th string in Dna
    Motifs ← (Motif1, ..., Motift)
    if SCORE(Motifs) < SCORE(BestMotifs)
      BestMotifs ← Motifs
  return BestMotifs
```

apply **Laplace's Rule of Succession** to form Profile from *Motif*₁... *Motif*_{*i*-1}}

Greedy motif search - in action

With matrices

Dna
 tt**ACCT**taac
 g**ATGT**ctgtc
 acg**GCGT**tag
 cccta**ACGA**g
 cgtcag**AGGT**

		Motifs	ACCT							
	A:	1+1	0+1	0+1	0+1		2/5	1/5	1/5	1/5
COUNT(Motifs)	C:	0+1	1+1	1+1	0+1		1/5	2/5	2/5	1/5
	G:	0+1	0+1	0+1	0+1		1/5	1/5	1/5	1/5
	T:	0+1	0+1	0+1	1+1		1/5	1/5	1/5	2/5
						PROFILE(Motifs)				

If we use profile matrix to compute the probabilities of all 4-mers in the second string of DNA

g ATG	ATGT	TGT c	GT ct	T ctg	ctgt	tgtc
$1/5^4$	$4/5^4$	$1/5^4$	$4/5^4$	$2/5^4$	$2/5^4$	$1/5^4$

Motifs
ACCT
ATGT

We get lucky and we choose the implanted 4-mers ATGT

	A:	2+1	0+1	0+1	0+1		3/6	1/6	1/6	1/6
COUNT(Motifs)	C:	0+1	1+1	1+1	0+1		1/6	2/6	2/6	1/6
	G:	0+1	0+1	1+1	0+1		1/6	1/6	2/6	1/6
	T:	0+1	1+1	0+1	2+1		1/6	2/6	1/6	3/6
						PROFILE(Motifs)				

Greedy motif search - in action

Dna
 tt**ACCT**taac
 g**ATGT**ctgtc
 acg**GCGT**tag
 cccta**ACGA**g
 cgtcag**AGGT**

acg**G** cg**GC** g**GCG** **GCGT** **CGT**t **GT**ta **T**tag
 $12/6^4$ $2/6^4$ $2/6^4$ $12/6^4$ $3/6^4$ $2/6^4$ $2/6^4$

This time we will assume that **ACGG** is selected instead the **GCGT**

		ACCT		
	<i>Motifs</i>	ATGT		
		acg G		
	A: 3+1 0+1 0+1 1+1		4/7 1/7 1/7 1/7	
<i>COUNT(Motifs)</i>	C: 0+1 2+1 1+1 0+1		1/7 3/7 2/7 1/7	
	G: 0+1 0+1 2+1 1+1	<i>PROFILE(Motifs)</i>	1/7 1/7 3/7 2/7	
	T: 0+1 1+1 0+1 2+1		1/7 2/7 1/7 3/7	

ccct ccta cta**A** ta**AC** a**ACG** **ACGA** **CGA**g
 $18/7^4$ $3/7^4$ $2/7^4$ $1/7^4$ $16/7^4$ $36/7^4$ $2/7^4$

While in this case the profile-most probable 4-mer is **ACGA**

		ACCT		
	<i>Motifs</i>	ATGT		
		acg G		
		ACGA		
	A: 4+1 0+1 0+1 0+1		5/8 1/8 1/8 2/8	
<i>COUNT(Motifs)</i>	C: 0+1 3+1 1+1 0+1		1/8 4/8 2/8 1/8	
	G: 0+1 0+1 3+1 1+1	<i>PROFILE(Motifs)</i>	1/8 1/8 4/8 2/8	
	T: 0+1 1+1 0+1 2+1		1/8 2/8 1/8 3/8	

Greedy motif search - in action

Dna tt**ACCT**taac
g**ATGT**ctgtc
acg**GCGT**tag
cccta**ACGA**g
cgtcag**AGGT**

cgtc	gtca	tcag	cag A	ag AG	g AGG	AGGT
$1/8^4$	$8/8^4$	$8/8^4$	$8/8^4$	$10/8^4$	$8/8^4$	$60/8^4$

Motifs **ACCT**
ATGT
acg**G**
ACGA
AGGT
CONSENSUS(*Motifs*) **ACGT**

Laplace's Rule of Succession has provided
a great improvement over the original
GreedyMotifSearch

Randomised Motif Search

Randomised algorithms may be nonintuitive because they lack the control of traditional algorithms

These algorithms are not guaranteed to return exact solutions, but they quickly find approximate solutions

Given a collection of strings DNA and an arbitrary $4 \times k$ matrix profile, we define $\text{MOTIF}(\text{Profile}, \text{DNA})$ as the collections of k -mers formed by the Profile-most probable k -mers in each sequence from DNA.

Considering profile and DNA

the profile-most 4-mer from each row of DNA produces the following 4-mers:

<i>Profile</i>	A: 4/5 0 0 1/5	<i>Dna</i>	ttaccttaac
	C: 0 3/5 1/5 0		gatgtctgtc
	G: 1/5 1/5 4/5 0		acggcgtag
	T: 0 1/5 0 4/5		ccctaacgag
			cgtcagaggt

MOTIFS(Profile, Dna)

tt acct taac
ga tg ctgtc
ac ggcg tag
cccta acg ag
cgtcag aggt

The **general idea** is that we can begin from a collection of randomly chosen k -mers Motif in DNA construct the profile and use this profile to generate a new collection of k -mers.

Randomised Motif Search

RANDOMIZEDMOTIFSEARCH(*Dna*, *k*, *t*)

randomly select *k*-mers *Motifs* = (*Motif*₁, ..., *Motif*_{*t*}) in each string from *Dna*

BestMotifs ← *Motifs*

while forever

Profile ← PROFILE(*Motifs*)

Motifs ← MOTIFS(*Profile*, *Dna*)

if SCORE(*Motifs*) < SCORE(*BestMotifs*)

BestMotifs ← *Motifs*

else

return *BestMotifs*

A single run may generate a poor set of motifs, then usually it is necessary run this algorithm thousand of times. In each run it begin from a new randomly selected set of *k*-mers and finally we select the best set of *k*-mers found in all these runs.

Randomised Motif Search - in action

Dna ttACCT**taac**
 gAT**GTct**gtc
 ccg**G**CGTtag
 c**acta**ACGAg
 cgtcag**AGGT**

We construct the profile matrix of the chosen 4-mers

<i>Motifs</i>				PROFILE(<i>Motifs</i>)				
t	a	a	c	A:	0.4	0.2	0.2	0.2
G	T	c	t	C:	0.2	0.4	0.2	0.2
c	c	g	G	G:	0.2	0.2	0.4	0.2
a	c	t	a	T:	0.2	0.2	0.2	0.4
A	G	G	T					

ttAC	tACC	ACCT	CCTt	CTta	Ttaa	taac
.0016	.0016	.0128	.0064	.0016	.0016	.0016
gATG	ATGT	TGTc	GTct	Tctg	ctgt	tgtc
.0016	.0128	.0016	.0032	.0032	.0032	.0016
ccgG	cgGC	gGCG	GCGT	CGTt	GTta	Ttag
.0064	.0036	.0016	.0128	.0032	.0016	.0016
cact	acta	ctaA	taAC	aACG	ACGA	CGAg
.0032	.0064	.0016	.0016	.0032	.0128	.0016
cgtc	gtca	tcag	cagA	agAG	gAGG	AGGT
.0016	.0016	.0016	.0032	.0032	.0032	.0128

We compute the probabilities of every 4-mer in DNA based on this profile matrix.

Dna tt**ACCT**taac
 g**ATGT**ctgtc
 ccg**GCGT**tag
 cacta**ACGA**g
 cgtcag**AGGT**

Randomised Motif Search

Starting from a **uniform distribution** as the first profile matrix is useless because no string is more probable than any other according to this profile and it does not provide any clues on what an implanted motif looks like.

Considering that if the strings in DNA were random, then the algorithm would start form a nearly uniform profile, and there would be nothing to work with. **The KEY observation is that the strings in DNA are NOT random.**

Gibbs Sampling

AlignACE Bioprospector

The randomise strategy can discarded all k-mers in each iteration, while Gibbs sampling is more conservative.

The algorithm starts from randomly chosen k-mers in each DNA, and select an integer between 1 and t and randomly changes a single k-mer.

```
ttaccttaac      ttaccttaac
gatatctgtc      gatatctgtc
acggcggttcg    acggcgttcg
ccctaaagag      ccctaaagag
cgtcagaggt      cgtcagaggt
```

RANDOMIZEDMOTIFSEARCH

(may change all *k*-mers in one step)

```
ttaccttaac      ttaccttaac
gatatctgtc      gatatctgtc
acggcggttcg    acggcggttcg
ccctaaagag      ccctaaagag
cgtcagaggt      cgtcagaggt
```

GIBBSSAMPLER

(changes one *k*-mer in one step)

Gibbs Sampling

GIBBSAMPLER(*Dna*, *k*, *t*, *N*)

randomly select *k*-mers *Motifs* = (*Motif*₁, ..., *Motif*_{*t*}) in each string from *Dna*

BestMotifs ← *Motifs*

for *j* ← 1 to *N*

i ← RANDOM(*t*)

Profile ← profile matrix formed from all strings in *Motifs* except for *Motif*_{*i*}

*Motif*_{*i*} ← *Profile*-randomly generated *k*-mer in the *i*-th sequence

if SCORE(*Motifs*) < SCORE(*BestMotifs*)

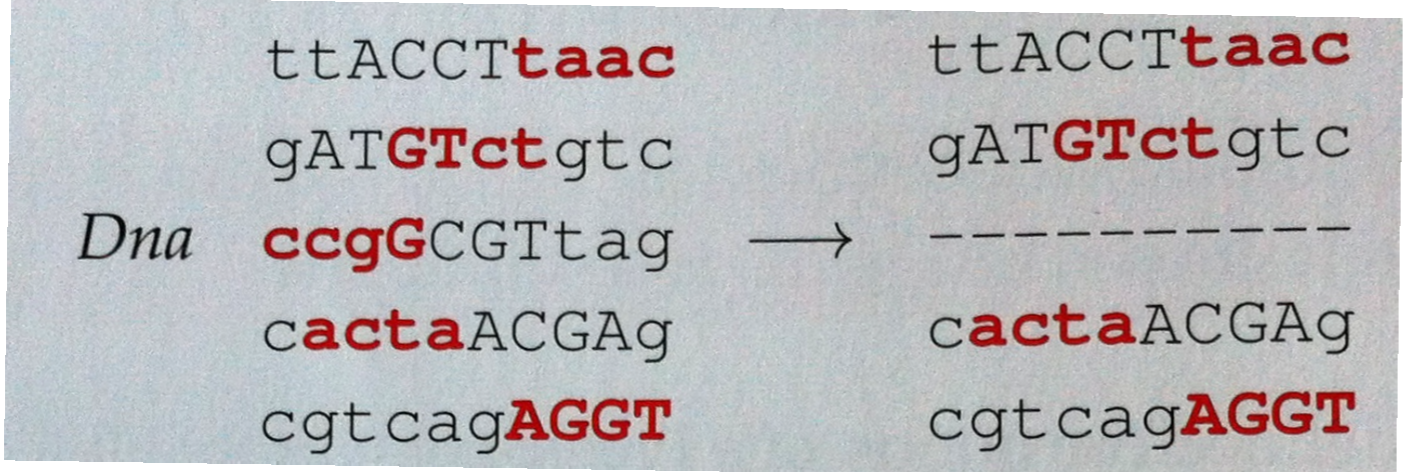
BestMotifs ← *Motifs*

return *BestMotifs*

Profile-randomly generated *k*-mer in the *i*-th sequence

Gibbs sampling uses an advice random number generator.

Gibbs Sampling - in action



Initial step, from a set of random k-mers, the algorithm select the third string for removal.

			t a a c
	<i>Motifs</i>		G T c t
			a c t a
			A G G T
	A: 2 1 1 1		A: 2/4 1/4 1/4 1/4
	C: 0 1 1 1		C: 0 1/4 1/4 1/4
<i>COUNT(Motifs)</i>	G: 1 1 1 0	<i>PROFILE(Motifs)</i>	G: 1/4 1/4 1/4 0
	T: 1 1 1 2		T: 1/4 1/4 1/4 2/4

ccgG	cgGC	gGCG	GCGT	CGTt	GTta	Ttag
0	0	0	1/128	0	1/256	0

Gibbs Sampling - in action

COUNT(<i>Motifs</i>)	A: 3 2 2 2	PROFILE(<i>Motifs</i>)	A: 3/8 2/8 2/8 2/8
	C: 1 2 2 2		C: 1/8 2/8 2/8 2/8
	G: 2 2 2 1		G: 2/8 2/8 2/8 1/8
	T: 2 2 2 3		T: 2/8 2/8 2/8 3/8

application of Laplace's Rule to the count matrix

ccgG	cgGC	gGCG	GCGT	CGTt	GTta	Ttag
4/8 ⁴	8/8 ⁴	8/8 ⁴	24/8 ⁴	12/8 ⁴	16/8 ⁴	8/8 ⁴

Gibbs sampling is not deterministic, then create a seven-face dice in which the total sum is equal to 80/8⁴

$$\text{RANDOM} \left(\frac{4/8^4}{80/8^4}, \frac{8/8^4}{80/8^4}, \frac{8/8^4}{80/8^4}, \frac{24/8^4}{80/8^4}, \frac{12/8^4}{80/8^4}, \frac{16/8^4}{80/8^4}, \frac{8/8^4}{80/8^4} \right)$$

$$= \text{RANDOM} \left(\frac{4}{80}, \frac{8}{80}, \frac{8}{80}, \frac{24}{80}, \frac{12}{80}, \frac{16}{80}, \frac{8}{80} \right)$$

The deleted string is now added but instead to use the **ccgG** mer, we roll the die and we obtain the mer **GCGT**

Gibbs Sampling - in action

Dna ttACCT**taac** -----
 gAT**GTct**gtc gAT**GTct**gtc
 ccg**GCGT**tag → ccg**GCGT**tag
 c**acta**ACGAg c**acta**ACGAg
 cgtcag**AGGT** cgtcag**AGGT**

Randomly selection of the deletion of the first DNA string

<i>Motifs</i>	G T c t	PROFILE(<i>Motifs</i>)	A: 2/4 0 0 1/4
	G C G T		C: 0 2/4 1/4 0
	a c t a		G: 2/4 1/4 2/4 0
	A G G T		T: 0 1/4 1/4 3/4

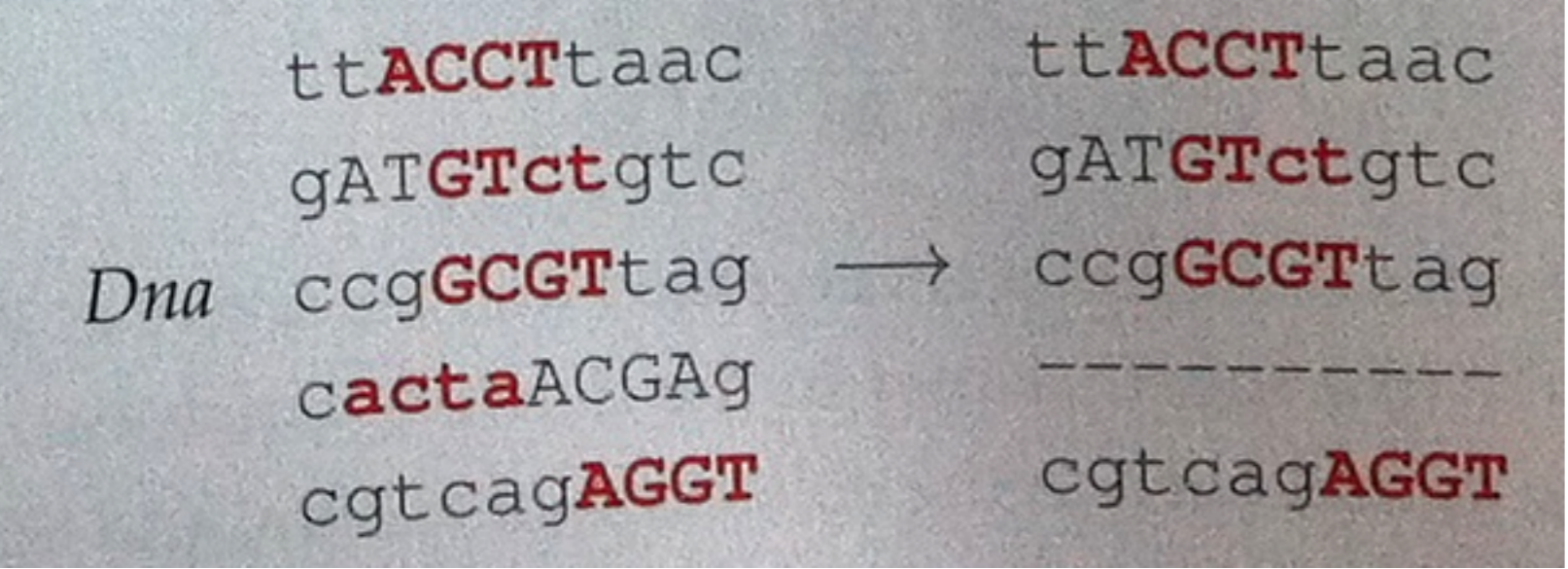
application of Laplace's Rule to the count matrix

COUNT(<i>Motifs</i>)	A: 3 1 1 2	PROFILE(<i>Motifs</i>)	A: 3/8 1/8 1/8 2/8
	C: 1 3 2 1		C: 1/8 3/8 2/8 1/8
	G: 3 2 3 1		G: 3/8 2/8 3/8 1/8
	T: 1 2 2 4		T: 1/8 2/8 2/8 4/8

ttAC	tACC	ACCT	CCTt	CTta	Ttaa	taac
$2/8^4$	$2/8^4$	$72/8^4$	$24/8^4$	$8/8^4$	$4/8^4$	$1/8^4$



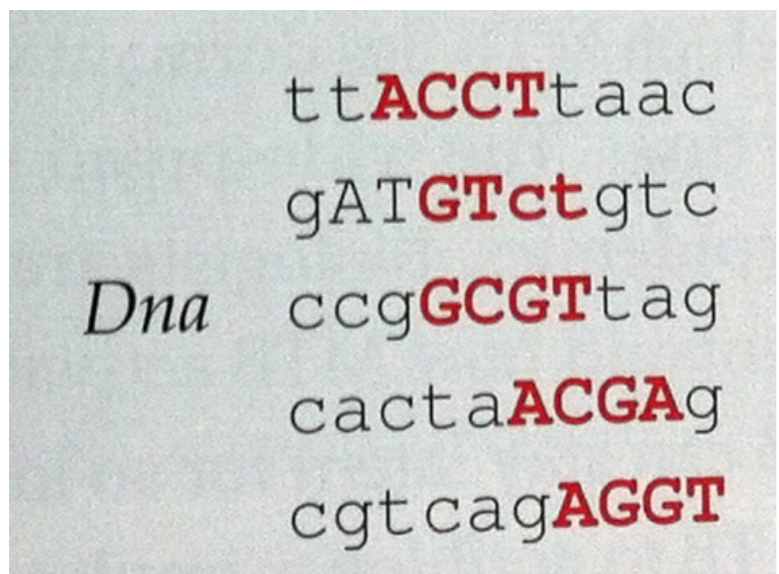
Gibbs Sampling - in action



Randomly selection of the deletion of the fourth DNA string

Motifs	A C C T				cact	acta	ctaA	taAC	aACG	ACGA	CGA	
		G	T	c	t	15/8 ⁴	9/8 ⁴	2/8 ⁴	1/8 ⁴	9/8 ⁴	27/8 ⁴	2/8 ⁴
	G	C	G	T								
	A	G	G	T								
COUNT(Motifs)	A: 3	1	1	1	PROFILE(Motifs)	A: 3/8	1/8	1/8	1/8			
	C: 1	3	3	1		C: 1/8	3/8	3/8	1/8			
	G: 3	2	3	1		G: 3/8	2/8	3/8	1/8			
	T: 1	2	1	5		T: 1/8	2/8	1/8	5/8			

Profile-randomly generated k-mer

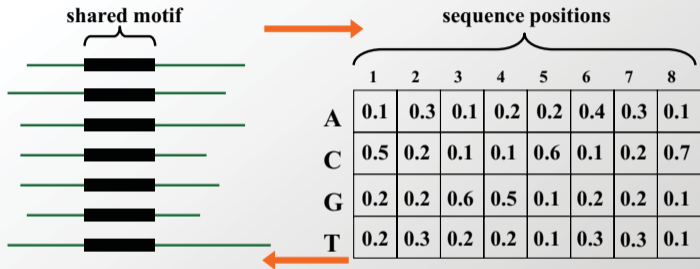


The algorithm starts to converge, but in some cases the algorithm can converge in a suboptimal solution.

A **local optimum** is a solution that is optimal within a small neighbouring set of solutions, which is in contact to the **global optimum** or the optimal solutions among all possible solutions.

Starting positions \Leftrightarrow Motif matrix

- given aligned sequences \rightarrow easy to compute profile matrix



- easy to find starting position probabilities \leftarrow given profile matrix

Key idea: Iterative procedure for estimating both, given uncertainty
(learning problem with hidden variables: the starting positions)

Basic Iterative Approach

Given: length parameter W , training set of sequences

set initial values for **motif**

do

→ re-estimate *starting-positions* from *motif*

→ re-estimate *motif* from *starting-positions*

until convergence (change $< \epsilon$)

return: *motif*, *starting-positions*

Representing Motif $M(k,c)$ and Background $B(c)$

- Assume motif has fixed width, W
- Motif represented by matrix of probabilities: $M(k,c)$
the probability of character c in column k

	1	2	3		
$M =$	A	0.1	0.5	0.2	
	C	0.4	0.2	0.1	(~CAG)
	G	0.3	0.1	0.6	
	T	0.2	0.2	0.1	

- Background represented by $B(c)$, frequency of each base

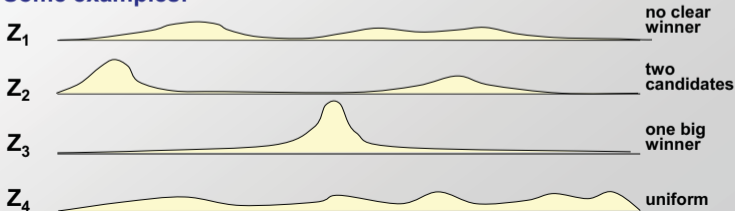
$B =$	A	0.26	
	C	0.24	(near uniform)
	G	0.23	(see also: di-nucleotide etc)
	T	0.27	

Representing the starting position probabilities (Z_{ij})

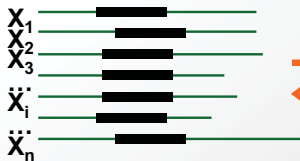
- the element Z_{ij} of the matrix Z represents the probability that the motif starts in position j in sequence i

	1	2	3	4
$Z =$ 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
seq4	0.1	0.5	0.1	0.3

Some examples:



Starting positions (Z_{ij}) \Leftrightarrow Motif matrix $M(k,c)$



Starting positions: Z_{ij}

M-step \rightarrow

\leftarrow E-step

	k=1	k=2	k=3	k=4	k=5	k=6	k=7	k=8
c=A	0.1	0.3	0.1	0.2	0.2	0.4	0.3	0.1
c=C	0.5	0.2	0.1	0.1	0.6	0.1	0.2	0.7
c=G	0.2	0.2	0.6	0.5	0.1	0.2	0.2	0.1
c=T	0.2	0.3	0.2	0.2	0.1	0.3	0.3	0.1

Motif: $M(k,c)$

- Z_{ij} : Probability that on sequence i , motif start at position j
- $M(k,c)$: Probability that k^{th} character of motif is letter c

- **Computing Z_{ij} matrix from $M(k,c)$ is straightforward**

- At each position, evaluate start probability by multiplying across the matrix

- **Three variations for re-computing motif $M(k,c)$ from Z_{ij} matrix**

- Expectation maximization \rightarrow All starts weighted by Z_{ij} prob distribution
- Gibbs sampling \rightarrow Single start for each seq X_i by sampling Z_{ij}
- Greedy approach \rightarrow Best start for each seq X_i by maximum Z_{ij}

E-step: Estimate Z_{ij} positions from matrix



Starting positions: Z_{ij}

	k=1	k=2	k=3	k=4	k=5	k=6	k=7	k=8
c=A	0.1	0.3	0.1	0.2	0.2	0.4	0.3	0.1
c=C	0.5	0.2	0.1	0.1	0.6	0.1	0.2	0.7
c=G	0.2	0.2	0.6	0.5	0.1	0.2	0.2	0.1
c=T	0.2	0.3	0.2	0.2	0.1	0.3	0.3	0.1

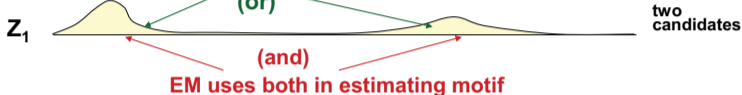
Motif: $M(k,c)$

Three examples for Greedy, Gibbs Sampling, EM

Greedy always picks maximum

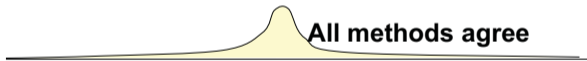
Gibbs sampling picks one at random
(or)

EM uses both in estimating motif



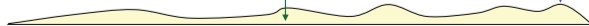
Z_2 All methods agree

one big winner



Greedy ignores most of the probability
Gibbs sampling rapidly converges to some choice

Z_3 uniform
EM averages over the entire sequence (slow/no convergence)



Calculating $P(X_i)$ when motif position is known

- Probability of training sequence X_i , given hypothesized start position j

$$\Pr(X_i | Z_{ij} = 1, M, B) = \underbrace{\prod_{k=1}^{j-1} B(X_{i,k})}_{\text{before motif}} \underbrace{\prod_{k=j}^{j+W-1} M(k-j+1, X_{i,k})}_{\text{motif}} \underbrace{\prod_{k=j+W}^L B(X_{i,k})}_{\text{after motif}}$$

- Example:

$$X_i = \text{G C } \boxed{\text{T G T}} \text{ A G} \quad B = \begin{array}{l} \text{A} \ 0.25 \\ \text{C} \ 0.25 \\ \text{G} \ 0.25 \\ \text{T} \ 0.25 \end{array} \quad M = \begin{array}{l} \text{A} \ 0.1 \ 0.5 \ 0.2 \\ \text{C} \ 0.4 \ 0.2 \ 0.1 \\ \text{G} \ 0.3 \ \boxed{0.1} \ 0.6 \\ \text{T} \ \boxed{0.2} \ 0.2 \ \boxed{0.1} \end{array}$$

$$\Pr(X_i | Z_{i3} = 1, M, B) =$$

$$B(\text{G}) \times B(\text{C}) \times M(1, \text{T}) \times M(2, \text{G}) \times M(3, \text{T}) \times B(\text{A}) \times B(\text{G}) =$$

$$0.25 \times 0.25 \times \boxed{0.2 \times 0.1 \times 0.1} \times 0.25 \times 0.25$$

Calculating the Z vector (using M)

- To estimate the starting positions in Z at step t

$$Z_{ij}^{(t)} = \Pr(Z_{ij} = 1 | X_i, M^{(t)}) = \frac{\overset{\text{likelihood}}{\Pr(X_i | Z_{ij} = 1, M^{(t)})} \overset{\text{prior}}{\Pr(Z_{ij} = 1)}}{\underset{\text{evidence}}{\Pr(X_i)}} \quad \text{(Bayes' rule)}$$

posterior

- At iteration t, calculate $Z_{ij}^{(t)}$ based on $M^{(t)}$
 - We just saw how to calculate $\Pr(X_i | Z_{ij}=1, M^{(t)})$
 - To obtain total probability $\Pr(X_i)$, sum over all starting positions

$$Z_{ij}^{(t)} = \frac{\Pr(X_i | Z_{ij} = 1, M^{(t)}) \cancel{\Pr(Z_{ij} = 1)}}{\sum_{k=1}^{L-W+1} \Pr(X_i | Z_{ik} = 1, M^{(t)}) \cancel{\Pr(Z_{ik} = 1)}}$$

- Assume uniform priors (motif eq likely to start at any position)

Calculating the Z vector: Example

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

	0	1	2	3
A	0.25	0.1	0.5	0.2
C	0.25	$\boxed{0.4}$	$\boxed{0.2}$	0.1
G	0.25	$\boxed{0.3}$	0.1	$\boxed{0.6}$
T	0.25	0.2	$\boxed{0.2}$	$\boxed{0.1}$

$$Z_{i1} = \boxed{0.3 \times 0.2 \times 0.1} \times 0.25 \times 0.25 \times 0.25 \times 0.25$$

$$Z_{i2} = 0.25 \times \boxed{0.4 \times 0.2 \times 0.6} \times 0.25 \times 0.25 \times 0.25$$

⋮

- then normalize so that

$$\sum_{j=1}^{L-W+1} Z_{ij} = 1$$

M-step: Max-likelihood motif from Z_{ij} positions



	k=1	k=2	k=3	k=4	k=5	k=6	k=7	k=8
e=A	0.1	0.3	0.1	0.2	0.2	0.4	0.3	0.1
e=C	0.5	0.2	0.1	0.1	0.6	0.1	0.2	0.7
e=G	0.2	0.2	0.6	0.5	0.1	0.2	0.2	0.1
e=T	0.2	0.3	0.2	0.2	0.1	0.3	0.3	0.1

Starting positions: Z_{ij}

Motif: $M(k,c)$

M-step example: Estimating $M(k,c)$ from Z_{ij}

$$X_1 = \mathbf{A} \quad \boxed{\mathbf{C} \ \mathbf{A} \ \mathbf{G}} \quad \mathbf{C} \ \mathbf{A}$$

$$Z_1 = 0.1 \quad 0.7 \quad 0.1 \quad 0.1$$

$$X_2 = \mathbf{A} \ \mathbf{G} \ \mathbf{G} \quad \boxed{\mathbf{C} \ \mathbf{A} \ \mathbf{G}}$$

$$Z_2 = 0.4 \quad 0.1 \quad 0.1 \quad 0.4$$

$$X_3 = \mathbf{T} \quad \boxed{\mathbf{C} \ \mathbf{A} \ \mathbf{G}} \quad \mathbf{T} \ \mathbf{C}$$

$$Z_3 = 0.2 \quad 0.6 \quad 0.1 \quad 0.1$$

$$M(1, A) = \frac{Z_{1,1} + Z_{1,3} + Z_{2,1} + Z_{3,3} + 1}{Z_{1,1} + Z_{1,2} + \dots + Z_{3,3} + Z_{3,4} + 4}$$

Em approach: Avg'em all
Gibbs sampling: Sample one
Greedy: Select max

- EM: sum over full probability

- $n_{1,A} = 0.1 + 0.1 + 0.4 + 0.1 = 0.7$

- $n_{1,C} = 0.7 + 0.4 + 0.6 = 1.7$

- $n_{1,G} = 0.1 + 0.1 + 0.1 + 0.1 = 0.4$

- $n_{1,T} = 0.2 = 0.2$

- Total: $T = 0.7 + 1.7 + 0.4 + 0.2 = 3.0$

- Normalize and add pseudo-counts

- $M(1,A) = (0.7+1)/(T+4) = 1.7/7 = 0.24$

- $M(1,C) = (1.7+1)/(T+4) = 2.7/7 = 0.39$

- $M(1,G) = (0.4+1)/(T+4) = 1.4/7 = 0.2$

- $M(1,T) = (0.2+1)/(T+4) = 1.2/7 = 0.17$

	1	2	3
A	0.24	0.39	0.21
C	0.39	0.21	0.18
G	0.2	0.24	0.44
T	0.17	0.16	0.16

- $M(k,c) =$

The EM Algorithm

- EM converges to a local maximum in the likelihood of the data given the model:

$$\prod_i \Pr(X_i | M, B)$$

- **Deterministic iterations max direction of ascent**
- **Usually converges in a small number of iterations**
- **Sensitive to initial starting point (i.e. values in M)**

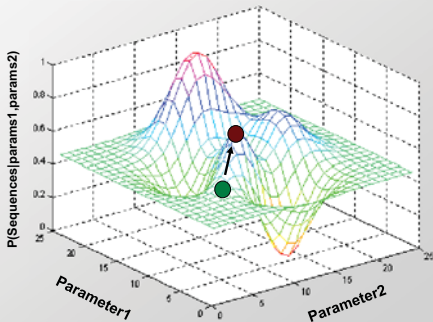
P(Seq|Model) Landscape

EM searches for parameters to increase $P(\text{seqs}|\text{parameters})$

Useful to think of $P(\text{seqs}|\text{parameters})$ as a **function of parameters**

EM starts at an **initial** set of parameters ●

And then “climbs uphill” until it reaches a **local maximum** ●



Where EM starts can make a big difference