

Algoritmos Aleatorios y Motivos Regulatorios

Bioinformática 2025-2

Universidad de Sonora

15 de octubre de 2025

Distintos motivos. TCGGGGATTTC

1	T	C	G	G	G	G	g	T	T	T	t	t
2	c	C	G	G	t	G	A	c	T	T	a	C
3	a	C	G	G	G	G	A	T	T	T	t	C
4	T	t	G	G	G	G	A	c	T	T	t	t
5	a	a	G	G	G	G	A	c	T	T	C	C
6	T	t	G	G	G	G	A	c	T	T	C	C
7	T	C	G	G	G	G	A	T	T	c	a	t
8	T	C	G	G	G	G	A	T	T	c	C	t
9	T	a	G	G	G	G	A	a	c	T	a	C
10	T	C	G	G	G	t	A	T	a	a	C	C

Distintos motivos

```
1 "atgacgggatactgataaaaaaagggggggggcgtacacattagataaacgtatgaagtacgttagactcggcgccgccg"  
2 "accctatTTTTTgagcagatttagtgacctggaaaaaaatttgagtacaaaactTTTccgaataaaaaaaggggggga"  
3 "tgagtatccctgggatgacttaaaaaaaggggggggtgctctcccattTTTgaatatgtaggatcattgccaggggtccga"  
4 "gtgagaattggatgaaaaaaaggggggggtccacgcaatcggaaccaacgcggacccaagggaagaccgataaaggaga"  
5 "tccctTTTgcgtaattgtgccgggaggctggttacgtaggggaagccctaacgggacttaaaaaaagggggggcttatag"  
6 "gtcaatcatgttcttTgtaatggatttaaaaaaaggggggggacgcgttggcgcacccaattcagtgTgggcgagcgcaa"  
7 "cggtTTTggcccttTtagaggcccccgtaaaaaaaggggggggcaattatgagagagctaattctatcgcgTgcgtTtcat"  
8 "aacttgagTTaaaaaaagggggggctggggcacatacaagaggagtcttcttatcagTTaatgctgtatgacactatgta"  
9 "ttggcccatTggctaaaagcccaacttgacaaatggaagatagaatccttgcataaaaaaagggggggaccgaaagggaag"  
10 "ctggTgagcaacgcagacattcttacgtgcatttagctcgcttcggggatcTaatagcacgaagcttaaaaaaaggggggga"
```

Distintos motivos

```
1 "atgaccgggatactgatAAAAAAGGGGGGggcgtacacattagataaacgtatgaagtacgttagactcggcgccgcg"
2 "accctatTTTTTgagcagatttagtgacctggaaaaaaatttgagtacaaaactTTTccgaataAAAAAAGGGGGGa"
3 "tgagtatccctgggatgacttAAAAAAGGGGGGtgctctcccgattTTTgaatatgtaggatcattcggcaggggccga"
4 "gtgagaattggatgAAAAAAGGGGGGtccacgcaatcggaaccaacgcgacccaaaggcaagaccgataaaggaga"
5 "tccTTTTTcggtaatgtgccgggaggctggttacgtaggggaagccctaacggacttaataAAAAAAGGGGGGcttatag"
6 "gtcaatcatgttcttgaatggatttAAAAAAGGGGGGgaccgcttggcgacccaaattcagtgtggcgagcgcaa"
7 "cggttttggccctttagaggccccgtAAAAAAGGGGGGcaattatgagagagctaatttatcgctgctgtttcat"
8 "aacttgagttAAAAAAGGGGGGctggggcacatacaagaggagtcttccttatcagttaatgctgtatgacactatgta"
9 "ttggcccatggcctaaagcccaacttgacaaatggaagatagaatccttgcatAAAAAAGGGGGGaccgaaagggaag"
10 "ctggtgagcaacgcagagattcttacgtgcattagctcgcttcggggatctaatagcacgaagcttAAAAAAGGGGGGa"
```

Distintos motivos

```
1 "atgaccgggatactgatAgAAgAAAGGttGGGggcggtacacattagataaacgtatgaagtacgttagactcggcgccgccg"
2 "accctatTTTTTgagcagatttagtgacctggaaaaaaatttgagtacaaaactTTTccgaatacAAtAAAcGGcGGGa"
3 "tgagtatccctgggatgacttAAAAtAAtGGaGtGGtgctctcccgatttttgaatatgtaggatcattcgccaggggccga"
4 "gtgagaattggatgcAAAAAAGGGattGtccacgcaatcggaaccaacgcggacccaaaggcaagaccgataaaggaga"
5 "tccTTTTTcggtaatgtgccgggaggctggttacgtaggggaagccctaacggacttaataAAtAAAGGaaGGGcttatag"
6 "gtcaatcatgttcttgtgaatggatttAAcAAtAAGGGctGGgaccgcttggcgacccaaattcagtgtggcgagcgcaa"
7 "cggTTTTTggccttggtagaggccccgtAtAAAcAAGGaGGGccaattatgagagagctaatttatcgctgctgtttcat"
8 "aacttgagttAAAAAtAGGGaGccctggggcacatacaagaggagtcttcttatcagttaatgctgtatgacactatgta"
9 "ttggcccatggcctaaaagcccaacttgacaaatggaagatagaatccttgcatActAAAAAGGaGcGGaccgaaagggaag"
10 "ctggtagcaacgcagagattcttacgtgcattagctcgcttcggggatctaatagcacgaagcttActAAAAAGGaGcGGa"
```

Nueva definición

Definición. Dada una colección de cadenas ADN y $d, k \in \mathbb{N}$, un k -mero es un (k, d) -motivo si aparece en cada cadena de ADN con a lo más d diferencias.

Motif Enumeration

```
MotifEnumeration(Dna, k, d)  
  Patterns  $\leftarrow$  an empty set  
  for each k-mer Pattern in the first string in Dna  
    for each k-mer Pattern' differing from Pattern by at most d mismatches  
      if Pattern' appears in each string from Dna with at most d mismatches  
        add Pattern' to Patterns  
  remove duplicates from Patterns  
  return Patterns
```

Matriz de motivos

Motifs

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

Matriz de motivos

Motifs	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											

Perfil de motivos

Motifs	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												
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
Profile(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	1	1	.9	.9	.1	0	0	0	0	0
	T:	.7	.2	0	0	.1	.1	0	.5	.8	.7	.3	.4

Perfil de motivos

Motifs	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												
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
Profile(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	1	1	.9	.9	.1	0	0	0	0	0
	T:	.7	.2	0	0	.1	.1	0	.5	.8	.7	.3	.4
Consensus(Motifs)	T	C	G	G	G	G	A	T	T	T	C	C	

Entropía y log de motivos

A cada columna de la matriz de motivos le corresponde una distribución de probabilidad, esto es, una colección de números no negativos cuya suma es 1. Una forma de medir la dispersión o incertidumbre de una distribución de probabilidad (p_1, \dots, p_N) es con la entropía.

$$H(p_1, \dots, p_n) = - \sum_{i=1}^N p_i \log_2 p_i$$

Entropía y logo de motivos

Motifs	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												
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
Profile(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	1	1	.9	.9	.1	0	0	0	0	0
	T:	.7	.2	0	0	.1	.1	0	.5	.8	.7	.3	.4
Consensus(Motifs)	T	C	G	G	G	G	A	T	T	T	C	C	



Nuevo problema

Problema. Dada una colección de cadenas, encontrar un conjunto de k -meros, uno de cada cadena, que minimice el puntaje de la matriz de motivos.

Un algoritmo de fuerza bruta requeriría checar todos los k meros posibles de cada cadena, esto es, $(n - k + 1)^t$ posibilidades, y para calcular el puntaje, se necesitarían kt pasos, lo cual implica que necesitaríamos $(n - k + 1)^t kt$ pasos. Esto es, un algoritmo de búsqueda bruta es de complejidad $O(n^t kt)$.

Nuevo problema

Necesitamos reformular este problema:

Motivos \rightarrow *Consensus(Motifs)*

Consensus(Motifs) \rightarrow *Motifs*

Dada una colección de k -meros Motivos y un k -mero *Pattern*, definimos $d(\textit{Pattern}, \textit{Motifs})$ como la suma de las distancias de Hamming entre *Pattern* y cada Motivo.

Distancia de Pattern a Motifs

<i>Motifs</i>	T	C	G	G	G	G	g	T	T	T	t	t	3
	c	C	G	G	t	G	A	c	T	T	a	C	4
	a	C	G	G	G	G	A	T	T	T	t	C	2
	T	t	G	G	G	G	A	c	T	T	t	t	4
	a	a	G	G	G	G	A	c	T	T	C	C	3
	T	t	G	G	G	G	A	c	T	T	C	C	2
	T	C	G	G	G	G	A	T	T	c	a	t	3
	T	C	G	G	G	G	A	T	T	c	C	t	2
	T	a	G	G	G	G	A	a	c	T	a	C	4
	T	C	G	G	G	t	A	T	a	a	C	C	<u>+ 3</u>
SCORE(<i>Motifs</i>)	3 + 4 + 0 + 0 + 1 + 1 + 1 + 5 + 2 + 3 + 6 + 4 = 30												
CONSENSUS(<i>Motifs</i>)	T	C	G	G	G	G	A	T	T	T	C	C	

Distancia de Pattern a Motifs

$$d(\text{Pattern}, \text{Motifs}) = \sum_{i=0}^t \text{HammingDistance}(\text{Pattern}, \text{Motifs})$$

$$\text{Score}(\text{Motifs}) = d(\text{Consensus}(\text{Motifs}), \text{Motifs})$$

Distancia de Pattern a Motifs

<i>Dna</i>	ttaccttAAC	1
	gATAtctgtc	1
	ACGgcgttcg	2
	ccctAAAgag	0
	cgtcAGAggt	1

Algoritmo Median String

MedianString(*Dna*, *k*)

distance $\leftarrow \infty$

Median $\leftarrow ""$

for each *k*-mer *Pattern* from AA...AA to TT...TT

if *distance* > *d*(*Pattern*, *Dna*)

distance $\leftarrow d(\textit{Pattern}, \textit{Dna})$

Median $\leftarrow \textit{Pattern}$

return *Median*

Calcular probabilidades con el perfil de motivos

Profile

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

$$\Pr(\text{ACGGGGATTACC}, \text{Profile}) = .2 \cdot .6 \cdot 1 \cdot 1 \cdot .9 \cdot .9 \cdot .9 \cdot .5 \cdot .8 \cdot .1 \cdot .4 \cdot .6 \\ = 0.000839808$$

$$\Pr(\text{TCGGGGATTCC} \mid \text{Profile}) = 0.7 \cdot 0.6 \cdot 1.0 \cdot 1.0 \cdot 0.9 \cdot 0.9 \cdot 0.9 \cdot 0.5 \cdot 0.8 \cdot 0.7 \cdot 0.4 \cdot 0.6 \\ = 0.0205753$$

GreedyMotifSearch

GreedyMotifSearch(*Dna*, *k*, *t*)

BestMotifs \leftarrow motif matrix formed by first *k*-mers in each string from *Dna*

for each *k*-mer *Motif* in the first string from *Dna*

*Motif*₁ \leftarrow *Motif*

for *i* = 2 to *t*

form *Profile* from motifs *Motif*₁, ..., *Motif*_{*i* - 1}

*Motif*_{*i*} \leftarrow *Profile*-most probable *k*-mer in the *i*-th string in *Dna*

Motifs \leftarrow (*Motif*₁, ..., *Motif*_{*t*})

if *Score*(*Motifs*) < *Score*(*BestMotifs*)

BestMotifs \leftarrow *Motifs*

return *BestMotifs*

GreedyMotifSearch

GreedyMotifSearch(*Dna*, *k*, *t*)

BestMotifs \leftarrow motif matrix formed by first *k*-mers in each string from *Dna*

for each *k*-mer *Motif* in the first string from *Dna*

*Motif*₁ \leftarrow *Motif*

for *i* = 2 to *t*

form *Profile* from motifs *Motif*₁, ..., *Motif*_{*i* - 1}

*Motif*_{*i*} \leftarrow *Profile*-most probable *k*-mer in the *i*-th string in *Dna*

Motifs \leftarrow (*Motif*₁, ..., *Motif*_{*t*})

if *Score*(*Motifs*) < *Score*(*BestMotifs*)

BestMotifs \leftarrow *Motifs*

return *BestMotifs*

Algunos problemas con GreedyMotifSearch

ttACCTaac

gATGTctgtc

acgGCGTtag

ccctaACGAg

cgtcagAGGT

Algunos problemas con GreedyMotifSearch

A: 1 0 0 0

C: 0 1 1 0

G: 0 0 0 0

T: 0 0 0 1

Algunos problemas con GreedyMotifSearch

Profile	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

$$\text{Pr}(\text{"TCGTGGATTTC"}, \text{Profile}) = .7 \cdot .6 \cdot 1 \cdot .0 \cdot .9 \cdot .9 \cdot .9 \cdot .5 \cdot .8 \cdot .7 \cdot .4 \cdot .6 = 0$$

Algunos problemas con GreedyMotifSearch

					T A A C					
					G T C T					
					A C T A					
					A G G T					
Count(Motifs)	A:	2	1	1	1	Profile(Motifs)	2/4	1/4	1/4	1/4
	C:	0	1	1	1		0	1/4	1/4	1/4
	G:	1	1	1	0		1/4	1/4	1/4	0
	T:	1	1	1	2		1/4	1/4	1/4	2/4

Algunos problemas con GreedyMotifSearch

Count(Motifs)

A:	2+1	1+1	1+1	1+1
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

Profile(Motifs)

	3/8	2/8	2/8	2/8
	1/8	2/8	2/8	2/8
	2/8	2/8	2/8	1/8
	2/8	2/8	2/8	3/8

Suavizado de Laplace

GreedyMotifSearch(*Dna*, *k*, *t*)

form a set of *k*-mers *BestMotifs* by selecting 1st *k*-mers in each string from *Dna*

for each *k*-mer *Motif* in the first string from *Dna*

*Motif*₁ ← *Motif*

 for *i* = 2 to *t*

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

*Motif*_{*i*} ← *Profile*-most probable *k*-mer in the *i*-th string in *Dna*

Motifs ← (*Motif*₁, ..., *Motif*_{*t*})

 if *Score*(*Motifs*) < *Score*(*BestMotifs*)

BestMotifs ← *Motifs*

output *BestMotifs*

Algunos problemas con GreedyMotifSearch

ttACCTaac

gATGTctgtc

acgGCGTtag

ccctaACGAg

cgtcagAGGT

GreedyMotifSearch con suavizado de Laplace

		Motifs ACCT			
Count(Motifs)	A:	1+1	0+1	0+1	0+1
	C:	0+1	1+1	1+1	0+1
	G:	0+1	0+1	0+1	0+1
	T:	0+1	0+1	0+1	1+1
		Profile(Motifs)			
		2/5	1/5	1/5	1/5
		1/5	2/5	2/5	1/5
		1/5	1/5	1/5	1/5
		1/5	1/5	1/5	2/5

Suavizado de Laplace

g**ATG**
 $1/5^4$

ATGT
 $4/5^4$

TGTc
 $1/5^4$

GTct
 $4/5^4$

Tctg
 $2/5^4$

ctgt
 $2/5^4$

tgtc
 $1/5^4$

Suavizado de Laplace

		Motifs			
		ACCT			
		ATGT			
Count(Motifs)	A:	2+1	0+1	0+1	0+1
	C:	0+1	1+1	1+1	0+1
	G:	0+1	0+1	1+1	0+1
	T:	0+1	1+1	0+1	2+1
		Profile(Motifs)			
		3/6	1/6	1/6	1/6
		1/6	2/6	2/6	1/6
		1/6	1/6	2/6	1/6
		1/6	2/6	1/6	3/6

Suavizado de Laplace

acg**G**
 $12/6^4$

cg**GC**
 $2/6^4$

g**GCG**
 $2/6^4$

GCGT
 $12/6^4$

CGTt
 $3/6^4$

GTta
 $2/6^4$

Ttag
 $2/6^4$

Suavizado de Laplace

Motifs **ACCT**
 ATGT
 acg**G**

Count(Motifs)	A:	3+1	0+1	0+1	1+1	Profile(Motifs)	4/7	1/7	1/7	1/7
	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		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

Suavizado de Laplace

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

Suavizado de Laplace

		Motifs			
		ACCT			
		ATGT			
		acgG			
		ACGA			
Count(Motifs)	A:	4+1	0+1	0+1	0+1
	C:	0+1	3+1	1+1	0+1
	G:	0+1	0+1	3+1	1+1
	T:	0+1	1+1	0+1	2+1
		Profile(Motifs)			
		5/8	1/8	1/8	2/8
		1/8	4/8	2/8	1/8
		1/8	1/8	4/8	2/8
		1/8	2/8	1/8	3/8

Suavizado de Laplace

cgtc
 $1/8^4$

gtca
 $8/8^4$

tcag
 $8/8^4$

cag**A**
 $8/8^4$

ag**AG**
 $10/8^4$

g**AGG**
 $8/8^4$

AGGT
 $60/8^4$

Suavizado de Laplace

	ACCT
	ATGT
Motifs	acgG
	ACGA
	AGGT
Consensus(Motifs)	ACGT

Suavizado de Laplace

Profile

A:	4/5	0	0	1/5
C:	0	3/5	1/5	0
G:	1/5	1/5	4/5	0
T:	0	1/5	0	4/5

Dna

ttaccttaac
gatgtctgtc
acggcgttag
ccctaacgag
cgtcagaggt

Suavizado de Laplace

Motifs(Profile, Dna)

```
ttaccttaac  
gatgtctgtc  
acggcgttag  
ccctaacgag  
cgtcagaggt
```


RandomizedMotifSearch

```
RandomizedMotifSearch(Dna, k, t)  
    randomly select k-mers Motifs = (Motif1, ..., Motift) 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
```