## Analisi di algoritmi per il Motif Finding

Tommaso Papini Gabriele Bani tommaso.papini1@stud.unifi.it gabriele.bani@stud.unifi.it



11 Dicembre 2015



## Un po' di background

#### DNA:

- √ sequenza di nucleotidi
- √ 4 tipi di nucleotide: A, T, C, G
- √ I-mer: sottosequenza di DNA di lunghezza I

### Motifs

In biologia può essere necessario ricavare certe sequenze di DNA "nascoste"

- ✓ pattern di nucleotidi ripetuti (I-mer)
- ✓ utili a capire determinati comportamenti biologici
  - sequenze di attivazione di geni specifici

## Il problema del Motif Finding

Il problema del Motif Finding consiste nel ricavare un set di t l-mer da un insieme di t sequenze di DNA.

### Input

- ✓ DNA: matrice di nucleotidi  $t \times n$ 
  - t sequenze di DNA
  - ognuna di lunghezza n
- √ I: lunghezza del motif cercato

### Output

 $\checkmark$   $s=(s_1,s_2,\ldots,s_t)$ : lista di t posizioni iniziali di l-mer il più simili tra loro

### Un primo esempio

CGGGGCTATGCAACTGGGTCGTCACATTCCCCTTTCGATA
TTTGAGGGTGCCCAATAAATGCAACTCCAAAGCGGACAAA
GGATGCAACTGATGCCGTTTGACGACCTAAATCAACGGCC
AAGGATGCAACTCCAGGAGCGCCTTTGCTGGTTCTACCTG
AATTTTCTAAAAAAGATTATAATGTCGGTCCATGCAACTTC
CTGCTGTACAACTGAGATCATGCTGCATGCAACTTTCAAC
TACATGATCTTTTGATGCAACTTGGATGAGGGAATGATGC

### Un primo esempio

CGGGGCTATGCAACTGGGTCGTCACATTCCCCTTTCGATA
TTTGAGGGTGCCCAATAAATGCAACTCCAAAGCGGACAAA
GGATGCAACTGATGCCGTTTGACGACCTAAATCAACGGCC
AAGGATGCAACTCCAGGAGCGCCTTTGCTGGTTCTACCTG
AATTTTCTAAAAAAGATTATAATGTCGGTCCATGCAACTTC
CTGCTGTACAACTGAGATCATGCTGCATGCAACTTTCAAC
TACATGATCTTTTGATGCAACTTGGATGAGGGAATGATGC

### Mutazioni random

CGGGGCTATcCAgCTGGGTCGTCACATTCCCCTTTCGATA
TTTGAGGGTGCCCAATAAggGCAACTCCAAAGCGGACAAA
GGATGgAtCTGATGCCGTTTGACGACCTAAATCAACGGCC
AAGGAaGCAACcCCAGGAGCGCCTTTGCTGGTTCTACCTG
AATTTTCTAAAAAGATTATAATGTCGGTCCtTGgAACTTC
CTGCTGTACAACTGAGATCATGCTGCATGCCAtTTTCAAC
TACATGATCTTTTGATGgcACTTGGATGAGGGAATGATGC

Come trovare l'I-mer più simile tra tutti?

### Allineamento

# CGGGGCT ATcCAgCT GGGTCGTCACATTCCCCTTTCGATA TTTGAGGGTGCCCAATAAggGCAACT CCAAAGCGGACAAA GGATGgAtCT GATGCCGTTTGACGACCTAAATCAACGGCC

AAGG<mark>AaGCAACc</mark>CCAGGAGCGCCTTTGCTGGTTCTACCTG

AATTTTCTAAAAAGATTATAATGTCGGTCC<mark>tTGgAACT</mark>TC CTGCTGTACAACTGAGATCATGCTGC<mark>ATGCcAtT</mark>TTCAAC

TACATGATCTTTTG<mark>ATGgcACT</mark>TGGATGAGGGAATGATGC

## Profilo e Consenso

		Α	Т	C	C	Α	G	C	T
		G	G	G	C	Α	Α	C	Т
Allineamento		Α	T	G	G	Α	T	C	/ T
		Α	Α	G	C	Α	Α	C	C
		T	Т	G	G	Α	Α	C	T
		Α	T	G	C	C	Α	T	Т
		Α	T	G	G	C	Α	C	T
Profilo	Α	5	1	0	0	5	5	0	0
	T	1	5	0	0	0	1	1	6
	G	1	1	6	3	0	1	0	0
	C	0	0	1	4	2	0	6	1
Consenso		Α	T	G	С	Α	Α	С	T

### Score

Come definire la "bontà" di un set di l-mer?

### Funzione score

Si definisce una funzione score sul vettore  $s = (s_1, s_2, \dots, s_t)$  di posizioni iniziali:

$$Score(s, DNA) = \sum_{j=1}^{I} M_{P(s)}(j)$$

dove

- $\checkmark$  P(s): matrice profile su s
- $\checkmark$   $M_{P(s)}(j)$ : elemento massimo nella colonna j-esima di P(s)

Si cerca il set di posizioni iniziali s che massimizzi Score(s, DNA)!

## Score: l'esempio di prima

		Α	Τ	C	C	Α	G	C	Т	
		G	G	G	C	Α	Α	C	T	
Allineamento		Α	T	G	G	Α	T	C	T	
		Α	Α	G	C	Α	Α	C	C	
		T	T	G	G	Α	Α	C	Т	
		Α	T	G	C	C	Α	T	Τ	
		Α	T	G	G	C	Α	C	T	
Profilo	Α	5	1	0	0	5	5	0	0	
	Т	1	5	0	0	0	1	1	6	
	G	1	1	6	3	0	1	0	0	
	C	0	0	1	4	2	0	6	1	
Consenso		Α	T	G	С	Α	Α	С	T	

$$Score(s, DNA) = 5 + 5 + 6 + 4 + 5 + 5 + 6 + 6 = 42$$

### Score

Quanto può valere lo score?

$$Score(s, DNA) = \begin{cases} I \cdot t, & \text{nel caso migliore} \\ \frac{I \cdot t}{4}, & \text{nel caso peggiore} \end{cases}$$

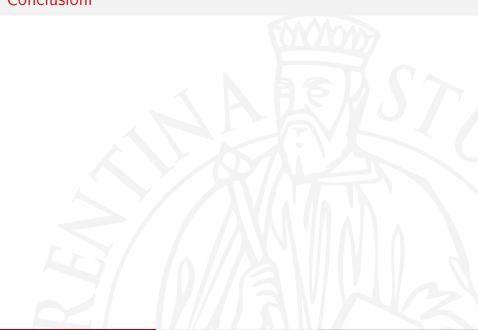
- ✓ It corrisponde al caso in cui tutti gli l-mer sono identici
- $\sqrt{\frac{lt}{4}}$  corrisponde al caso in cui gli l-mer siano diversi in tutte le posizioni

## Algoritmi brute force

# Algoritmi greedy

# Algoritmi randomizzati

## Conclusioni





Domande? Grazie!



Domande? Granie!



Domande? Grazie!