

Allineamento di Sequenze *(e strutture dati)*

Alberto Policriti

Dipartimento di Matematica e Informatica
Istituto di Genomica Applicata

7 Novembre, 2017

THE PLAN

BASICS QUESTIONS

- ① what does it mean to *align*?
- ② why should I care about *algorithms* to align?
- ③ data structures ... why?
- ④ time, space, correctness, completeness ... is this *all* bioinformaticians are interested on?

OUR (OPERATIVE) MEANING TO ALIGNMENT

ALIGNING MEANS ...

... converting

WHAT DO WE MEAN BY *aligning* TWO SEQUENCES?

DEFINITION

An alignment of σ_1 and σ_2 is a *string* on the alphabet $\{\mathbf{m}, \mathbf{i}, \mathbf{d}, \mathbf{s}\}$

HOW DO WE DEFINE AN *alignment*?

Forget about **how to find** an alignment or whether it is **optimal** (?)

EXAMPLE

$\sigma_1 \equiv$	a	c	g	t	c	a	t	c	a	
	i	m	s	m	m	s	d	m	m	m
$\sigma_2 \equiv$	t	a	a	g	t	g	t	c	a	

WHAT DO WE MEAN BY *aligning* TWO SEQUENCES?

DEFINITION

An alignment of σ_1 and σ_2 is a *string* on the alphabet $\{\mathbf{m}, \mathbf{i}, \mathbf{d}, \mathbf{s}\}$

HOW DO WE DEFINE AN *alignment*?

Forget about **how to find** an alignment or whether it is **optimal** (?)

EXAMPLE

$\sigma_1 \equiv$	a	c	g	t	c	a	t	c	a	
	i	m	s	m	m	s	d	m	m	m
$\sigma_2 \equiv$	t	a	a	g	t	g	t	c	a	

WHAT DO WE MEAN BY *aligning* TWO SEQUENCES?

DEFINITION

An alignment of σ_1 and σ_2 is a *string* on the alphabet $\{\mathbf{m}, \mathbf{i}, \mathbf{d}, \mathbf{s}\}$

HOW DO WE DEFINE AN *alignment*?

Forget about **how to find** an alignment or whether it is **optimal** (?)

EXAMPLE

$\sigma_1 \equiv$	a	c	g	t	c	a	t	c	a	
	i	m	s	m	m	s	d	m	m	m
$\sigma_2 \equiv$	t	a	a	g	t	g		t	c	a

WHAT DO WE MEAN BY *aligning* TWO SEQUENCES?HOW DO WE DEFINE AN *alignment*?Forget about **how to find** an alignment or whether it is **optimal** (?)

EXAMPLE

$\sigma_1 \equiv$	<i>a</i>	<i>c</i>	<i>g</i>	<i>t</i>	<i>c</i>	<i>a</i>	<i>t</i>	<i>c</i>	<i>a</i>
	i	m	s	m	m	s	d	m	m
$\sigma_2 \equiv$	<i>t</i>	<i>a</i>	<i>a</i>	<i>g</i>	<i>t</i>	<i>g</i>		<i>t</i>	<i>c</i>

AN ALIGNMENT IS A *string*Alignment \Leftrightarrow **edit-string**In fact is a “*program*” converting σ_1 into σ_2 (or viceversa).

TERMS AND MEANING

Ranking ALIGNMENTS: COST AND SCORE

We can associate a **cost** to an edit-string (then we want ↓)

We can associate a **score** to an edit-string (then we want ↑)

TERMS AND MEANING

Ranking ALIGNMENTS: COST AND SCORE

We can associate a **cost** to an edit-string (then we want ↓)

We can associate a **score** to an edit-string (then we want ↑)

ASSIGNING A SCORE

MATCHES (good): higher score (e.g. 1)

SUBSTITUTIONS (bad): lower score (e.g. less than 0)

INSERTION/DELETION (bad): lower score (e.g. less than 0)

TERMS AND MEANING

Ranking ALIGNMENTS: COST AND SCORE

We can associate a **cost** to an edit-string (then we want ↓)

We can associate a **score** to an edit-string (then we want ↑)

EXAMPLE

Purines and Pyrimidines are similar, hence a (non-trivial) **score matrix** is the following:

	<i>a</i>	<i>c</i>	<i>g</i>	<i>t</i>
<i>a</i>	2	-1	1	-1
<i>c</i>	-1	2	-1	1
<i>g</i>	1	-1	2	-1
<i>t</i>	-1	1	-1	2

TERMS AND MEANING

Ranking ALIGNMENTS: COST AND SCORE

We can associate a **cost** to an edit-string (then we want ↓)

We can associate a **score** to an edit-string (then we want ↑)

EXAMPLE

Affine gap penalty:

penalize the **opening** of a gap
and
penalise its **extension**

(... how exactly?)

TERMS AND MEANING

GLOBAL

The edit-string converts the *entire* σ_2 in *entire* σ_1 .

LOCAL

The edit-string converts the *entire* σ_2 in *a portion* of σ_1 .

TERMS AND MEANING

DISTANCE

Think of the *cost* as a *distance* among strings.

TERMS AND MEANING

DISTANCE

Think of the *cost* as a *distance* among strings.

EXAMPLE

Levensthein distance: cost 0 for **m**, 1 for the others.

$$\begin{array}{r}
 \sigma_1 \equiv \quad a \quad c \quad g \quad t \quad c \quad a \quad t \quad c \quad a \\
 \quad \quad \mathbf{i} \quad \mathbf{m} \quad \mathbf{s} \quad \mathbf{m} \quad \mathbf{m} \quad \mathbf{s} \quad \mathbf{d} \quad \mathbf{m} \quad \mathbf{m} \quad \mathbf{m} \Rightarrow d_L(\sigma_1, \sigma_2) = 4. \\
 \sigma_2 \equiv \quad t \quad a \quad a \quad g \quad t \quad g \quad \quad \quad t \quad c \quad a
 \end{array}$$

Hamming distance: **i**'s and **d**'s not allowed.

$$\begin{array}{r}
 \sigma_1 \equiv \quad a \quad c \quad g \quad t \quad c \quad a \quad t \quad c \quad a \\
 \quad \quad \mathbf{s} \quad \mathbf{s} \quad \mathbf{s} \quad \mathbf{s} \quad \mathbf{s} \quad \mathbf{s} \quad \mathbf{m} \quad \mathbf{m} \quad \mathbf{m} \Rightarrow d_H(\sigma_1, \sigma_2) = 6 \\
 \sigma_2 \equiv \quad t \quad a \quad a \quad g \quad t \quad g \quad t \quad c \quad a
 \end{array}$$

Hamming distance can be used on equal length strings only.

TERMS AND MEANING

Ranking ALIGNMENTS: COST AND SCORE

We can associate a **cost** to an edit-string (then we want ↓)

We can associate a **score** to an edit-string (then we want ↑)

GLOBAL

The edit-string converts the *entire* σ_2 in *entire* σ_1 .

LOCAL

The edit-string converts the *entire* σ_2 in *a portion* of σ_1 .

DISTANCE

Think of the **cost** as a *distance* among strings.

RULE OF THUMB

WHICH DISTANCE SHOULD I USE?

Long strings \Leftrightarrow Levenstein

Short strings \Leftrightarrow Hamming

RULE OF THUMB

WHICH DISTANCE & TOOL SHOULD I USE?

Long strings \Leftrightarrow Levenstein \Leftrightarrow S.W., BLAST, BWA-MEM, ...

Short strings \Leftrightarrow Hamming \Leftrightarrow Erne, BWA, bowtie, ...

RULE OF THUMB

WHICH DISTANCE & TOOL & COMPUTER SHOULD I USE?

Long strings \Leftrightarrow Levensthein \Leftrightarrow S.W., BLAST, BWA-MEM, ...

Short strings \Leftrightarrow Hamming \Leftrightarrow Erne, BWA, bowtie, ...

Always ... a powerful one (better if parallel)

$|\sigma|$ stands for the length of σ .

$O(|\sigma|)$ stands for *proportional* to the length of σ .

GLOBAL ALIGNMENT: COMPLEXITY

Best Hamming distance	Levensthein distance d	Levensthein
$O(\sigma_1 + \sigma_2)$	$O(\sigma_1 + \sigma_2)$	$O(\sigma_1 * \sigma_2)$

$|\sigma|$ stands for the length of σ .

$O(|\sigma|)$ stands for *proportional* to the length of σ .

GLOBAL ALIGNMENT: COMPLEXITY

Best Hamming distance	Levensthein distance d	Levensthein
$O(\sigma_1 + \sigma_2)$	$O(\sigma_1 + \sigma_2)$	$O(\sigma_1 * \sigma_2)$

Denote by ρ the reference (i.e. think of ρ as *large*).

LOCAL ALIGNMENT: COMPLEXITY

Best Hamming distance	Levensthein distance d	Levensthein
$O(\sigma_1 + \rho)$	$O(\sigma_1 + \rho)$	$O(\sigma_1 * \rho)$

PROGRAMMAZIONE DINAMICA PER L'ALLINEAMENTO
GLOBALE

IDEA!

Se dobbiamo trovare una stringa (i.e. η) determiniamone un carattere alla volta.

- Immagino di aver determinato $\eta[1, \dots, k]$ e voglio $\eta[k + 1]$. Diciamo che $\eta[1, \dots, k]$ allinea $\sigma_1[1, \dots, i - 1]$ e $\sigma_2[1, \dots, j - 1]$.
- 4 casi: $\eta[k + 1] \in \{\mathbf{m}, \mathbf{s}, \mathbf{i}, \mathbf{d}\}$.
- Se $\eta[k + 1] \in \{\mathbf{m}, \mathbf{s}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i]$ e $\sigma_2[1, \dots, j]$.
- Se $\eta[k + 1] \in \{\mathbf{i}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i - 1]$ e $\sigma_2[1, \dots, j]$.
- Se $\eta[k + 1] \in \{\mathbf{d}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i]$ e $\sigma_2[1, \dots, j - 1]$.

PROGRAMMAZIONE DINAMICA PER L'ALLINEAMENTO
GLOBALE

IDEA!

Se dobbiamo trovare una stringa (i.e. η) determiniamone un carattere alla volta.

- Immagino di aver determinato $\eta[1, \dots, k]$ e voglio $\eta[k + 1]$. Diciamo che $\eta[1, \dots, k]$ allinea $\sigma_1[1, \dots, i - 1]$ e $\sigma_2[1, \dots, j - 1]$.
- 4 casi: $\eta[k + 1] \in \{\mathbf{m}, \mathbf{s}, \mathbf{i}, \mathbf{d}\}$.
- Se $\eta[k + 1] \in \{\mathbf{m}, \mathbf{s}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i]$ e $\sigma_2[1, \dots, j]$.
- Se $\eta[k + 1] \in \{\mathbf{i}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i - 1]$ e $\sigma_2[1, \dots, j]$.
- Se $\eta[k + 1] \in \{\mathbf{d}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i]$ e $\sigma_2[1, \dots, j - 1]$.

PROGRAMMAZIONE DINAMICA PER L'ALLINEAMENTO
GLOBALE

IDEA!

Se dobbiamo trovare una stringa (i.e. η) determiniamone un carattere alla volta.

- Immagino di aver determinato $\eta[1, \dots, k]$ e voglio $\eta[k + 1]$. Diciamo che $\eta[1, \dots, k]$ allinei $\sigma_1[1, \dots, i - 1]$ e $\sigma_2[1, \dots, j - 1]$.
- 4 casi: $\eta[k + 1] \in \{\mathbf{m}, \mathbf{s}, \mathbf{i}, \mathbf{d}\}$.
- Se $\eta[k + 1] \in \{\mathbf{m}, \mathbf{s}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i]$ e $\sigma_2[1, \dots, j]$.
- Se $\eta[k + 1] \in \{\mathbf{i}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i - 1]$ e $\sigma_2[1, \dots, j]$.
- Se $\eta[k + 1] \in \{\mathbf{d}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i]$ e $\sigma_2[1, \dots, j - 1]$.

PROGRAMMAZIONE DINAMICA PER L'ALLINEAMENTO
GLOBALE

IDEA!

Se dobbiamo trovare una stringa (i.e. η) determiniamone un carattere alla volta.

- Immagino di aver determinato $\eta[1, \dots, k]$ e voglio $\eta[k + 1]$.
Diciamo che $\eta[1, \dots, k]$ allinei $\sigma_1[1, \dots, i - 1]$ e $\sigma_2[1, \dots, j - 1]$.
- 4 casi: $\eta[k + 1] \in \{\mathbf{m}, \mathbf{s}, \mathbf{i}, \mathbf{d}\}$.
- Se $\eta[k + 1] \in \{\mathbf{m}, \mathbf{s}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i]$ e $\sigma_2[1, \dots, j]$.
- Se $\eta[k + 1] \in \{\mathbf{i}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i - 1]$ e $\sigma_2[1, \dots, j]$.
- Se $\eta[k + 1] \in \{\mathbf{d}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i]$ e $\sigma_2[1, \dots, j - 1]$.

PROGRAMMAZIONE DINAMICA PER L'ALLINEAMENTO
GLOBALE

IDEA!

Se dobbiamo trovare una stringa (i.e. η) determiniamone un carattere alla volta.

- Immagino di aver determinato $\eta[1, \dots, k]$ e voglio $\eta[k + 1]$. Diciamo che $\eta[1, \dots, k]$ allinei $\sigma_1[1, \dots, i - 1]$ e $\sigma_2[1, \dots, j - 1]$.
- 4 casi: $\eta[k + 1] \in \{\mathbf{m}, \mathbf{s}, \mathbf{i}, \mathbf{d}\}$.
- Se $\eta[k + 1] \in \{\mathbf{m}, \mathbf{s}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i]$ e $\sigma_2[1, \dots, j]$.
- Se $\eta[k + 1] \in \{\mathbf{i}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i - 1]$ e $\sigma_2[1, \dots, j]$.
- Se $\eta[k + 1] \in \{\mathbf{d}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i]$ e $\sigma_2[1, \dots, j - 1]$.

PROGRAMMAZIONE DINAMICA PER L'ALLINEAMENTO
GLOBALE

IDEA!

Se dobbiamo trovare una stringa (i.e. η) determiniamone un carattere alla volta.

- Immagino di aver determinato $\eta[1, \dots, k]$ e voglio $\eta[k + 1]$. Diciamo che $\eta[1, \dots, k]$ allinea $\sigma_1[1, \dots, i - 1]$ e $\sigma_2[1, \dots, j - 1]$.
- 4 casi: $\eta[k + 1] \in \{\mathbf{m}, \mathbf{s}, \mathbf{i}, \mathbf{d}\}$.
- Se $\eta[k + 1] \in \{\mathbf{m}, \mathbf{s}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i]$ e $\sigma_2[1, \dots, j]$.
- Se $\eta[k + 1] \in \{\mathbf{i}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i - 1]$ e $\sigma_2[1, \dots, j]$.
- Se $\eta[k + 1] \in \{\mathbf{d}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i]$ e $\sigma_2[1, \dots, j - 1]$.

PROGRAMMAZIONE DINAMICA PER L'ALLINEAMENTO
GLOBALE

IDEA!

Se dobbiamo trovare una stringa (i.e. η) determiniamone un carattere alla volta.

- Immagino di aver determinato $\eta[1, \dots, k]$ e voglio $\eta[k + 1]$. Diciamo che $\eta[1, \dots, k]$ allinea $\sigma_1[1, \dots, i - 1]$ e $\sigma_2[1, \dots, j - 1]$.
- 4 casi: $\eta[k + 1] \in \{\mathbf{m}, \mathbf{s}, \mathbf{i}, \mathbf{d}\}$.
- Se $\eta[k + 1] \in \{\mathbf{m}, \mathbf{s}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i]$ e $\sigma_2[1, \dots, j]$.
- Se $\eta[k + 1] \in \{\mathbf{i}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i - 1]$ e $\sigma_2[1, \dots, j]$.
- Se $\eta[k + 1] \in \{\mathbf{d}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i]$ e $\sigma_2[1, \dots, j - 1]$.

PROGRAMMAZIONE DINAMICA PER L'ALLINEAMENTO
GLOBALE

ALGORITMO RICORSIVO

... esponenziale

(vedi anche Cabada A., Nieto J.J., Torres A.: *An exact formula for the number of alignments between two DNA sequences.*)

Algorithm 1: *allineamento_exp*(σ_1, σ_2, i, j)

```
if  $i == 0$  or  $j == 0$  then // casi base della ricorsione
  | return  $i + j$ 
if  $\sigma_1[i] == \sigma_2[j]$  then
  |  $c_m \leftarrow \text{allineamento\_exp}(\sigma_1, \sigma_2, i - 1, j - 1)$ ;
else
  |  $c_s \leftarrow \text{allineamento\_exp}(\sigma_1, \sigma_2, i - 1, j - 1) + 1$ ;
 $c_i \leftarrow \text{allineamento\_exp}(\sigma_1, \sigma_2, i - 1, j) + 1$ ;
 $c_d \leftarrow \text{allineamento\_exp}(\sigma_1, \sigma_2, i, j - 1) + 1$ ;
return  $\min(c_m, c_s, c_i, c_d)$ ;
```

PROGRAMMAZIONE DINAMICA PER L'ALLINEAMENTO
GLOBALE

(SUPER)OTTIMIZZAZIONE

Posso *tenere nota* dei risultati di tutti i possibili allineamenti parziali in una matrice.

$$A(i, j) = \text{(costo dell')} \text{ allineamento di } \sigma_1[1, \dots, i] \text{ con } \sigma_2[1, \dots, j]$$

PROGRAMMAZIONE DINAMICA PER L'ALLINEAMENTO
GLOBALE

(SUPER)OTTIMIZZAZIONE

Posso *tenere nota* dei risultati di tutti i possibili allineamenti parziali in una matrice.

$A(i, j)$ = (costo dell') allineamento di $\sigma_1[1, \dots, i]$ con $\sigma_2[1, \dots, j]$

ϵ	ϵ	j	
i		$d_L(i, j)$	

PROGRAMMAZIONE DINAMICA PER L'ALLINEAMENTO GLOBALE

EXAMPLE

Matrice di programmazione dinamica (da non confondere con la matrice dei costi di allineamento).

$$\begin{pmatrix}
 & \epsilon & a & c & g & t & c & a & t & c & a \\
 \epsilon & 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\
 t & 1 & & & & & & & & & \\
 a & 2 & & & & & & & & & \\
 a & 3 & & & & & & & & & \\
 g & 4 & & & & & & & & & \\
 t & 5 & & & & & & & & & \\
 g & 6 & & & & & & & & & \\
 t & 7 & & & & & & & & & \\
 c & 8 & & & & & & & & & \\
 a & 9 & & & & & & & & &
 \end{pmatrix}$$

PROGRAMMAZIONE DINAMICA PER L'ALLINEAMENTO
GLOBALE

EXAMPLE

Matrice di programmazione dinamica (da non confondere con la matrice dei costi di allineamento).

$$\left(\begin{array}{c|cccccccccc} & \epsilon & a & c & g & t & c & a & t & c & a \\ \hline \epsilon & 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ t & 1 & 1 & & & & & & & & \\ a & 2 & & & & & & & & & \\ a & 3 & & & & & & & & & \\ g & 4 & & & & & & & & & \\ t & 5 & & & & & & & & & \\ g & 6 & & & & & & & & & \\ t & 7 & & & & & & & & & \\ c & 8 & & & & & & & & & \\ a & 9 & & & & & & & & & \end{array} \right)$$

PROGRAMMAZIONE DINAMICA PER L'ALLINEAMENTO GLOBALE

EXAMPLE

Matrice di programmazione dinamica (da non confondere con la matrice dei costi di allineamento).

$$\begin{pmatrix}
 & \epsilon & a & c & g & t & c & a & t & c & a \\
 \epsilon & 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\
 t & 1 & 1 & & & & & & & & \\
 a & 2 & 1 & & & & & & & & \\
 a & 3 & & & & & & & & & \\
 g & 4 & & & & & & & & & \\
 t & 5 & & & & & & & & & \\
 g & 6 & & & & & & & & & \\
 t & 7 & & & & & & & & & \\
 c & 8 & & & & & & & & & \\
 a & 9 & & & & & & & & &
 \end{pmatrix}$$

PROGRAMMAZIONE DINAMICA PER L'ALLINEAMENTO GLOBALE

EXAMPLE

Matrice di programmazione dinamica (da non confondere con la matrice dei costi di allineamento).

$$\begin{pmatrix}
 & \epsilon & a & c & g & t & c & a & t & c & a \\
 \epsilon & 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\
 t & 1 & 1 & & & & & & & & \\
 a & 2 & 1 & & & & & & & & \\
 a & 3 & 2 & & & & & & & & \\
 g & 4 & & & & & & & & & \\
 t & 5 & & & & & & & & & \\
 g & 6 & & & & & & & & & \\
 t & 7 & & & & & & & & & \\
 c & 8 & & & & & & & & & \\
 a & 9 & & & & & & & & &
 \end{pmatrix}$$

PERCHÉ È COMPLESSO?

DOMANDA:

Cosa ho veramente computato?

ALLINEAMENTO LOCALE

DEFINITION

L'allineamento locale (a distanza $\leq d$) di σ_1 in σ_2 è l'insieme di tutte le coppie $\langle i, \eta_i \rangle$ tali che η_i è l'allineamento di σ_1 e di un *prefisso* di $\sigma_2[i, \dots]$ (a distanza $\leq d$).

N.B.

σ_1 *pattern (query)* σ_2 *testo (reference, data base, ...)* e ci interessano i **prefissi dei suffissi di σ_2** .

BASTA ...

ALLINEAMENTO LOCALE

DEFINITION

L'allineamento locale (a distanza $\leq d$) di σ_1 in σ_2 è l'insieme di tutte le coppie $\langle i, \eta_i \rangle$ tali che η_i è l'allineamento di σ_1 e di un *prefisso* di $\sigma_2[i, \dots]$ (a distanza $\leq d$).

N.B.

σ_1 *pattern (query)* σ_2 *testo (reference, data base, ...)* e ci interessano i **prefissi dei suffissi di σ_2** .

BASTA ...

... compilare la matrice di allineamento ponendo tutti 0 nella prima riga.

ALLINEAMENTO LOCALE

$$\left(\begin{array}{c|cccccccccc} & \epsilon & a & c & g & t & c & a & t & c & a \\ \hline \epsilon & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ t & 1 & 1 & 1 & 1 & 0 & & & & & \\ a & 2 & & & & & & & & & \\ a & 3 & & & & & & & & & \\ g & 4 & & & & & & & & & \\ t & 5 & & & & & & & & & \\ g & 6 & & & & & & & & & \\ t & 7 & & & & & & & & & \\ c & 8 & & & & & & & & & \\ a & 9 & & & & & & & & & \end{array} \right)$$

ALLINEAMENTO LOCALE

$$\left(\begin{array}{c|ccccccccc} & \epsilon & a & c & g & t & c & a & t & c & a \\ \hline \epsilon & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ t & 1 & 1 & 1 & 1 & 0 & & & & & \\ a & 2 & & & & & & & & & \\ a & 3 & & & & & & & & & \\ g & 4 & & & & & & & & & \\ t & 5 & & & & & & & & & \\ g & 6 & & & & & & & & & \\ t & 7 & & & & & & & & & \\ c & 8 & & & & & & & & & \\ a & 9 & & & & & & & & & \end{array} \right)$$

Abbiamo introdotto un *gap iniziale a costo zero*

Funzioni di costo per i gap

- Lineare: un gap di lunghezza n ha costo pari a $\gamma(n) \equiv n \cdot g$
- Convesso: un gap di lunghezza n ha costo pari a $\gamma(n)$ tale che

$$\gamma(i + 1) - \gamma(i) \leq \gamma(i) - \gamma(i - 1)$$

- Affine: un gap di lunghezza n ha costo pari a $\gamma(n) \equiv g_o + (n - 1) \cdot g_e$

GAP

Funzioni di costo per i gap

- Lineare: un gap di lunghezza n ha costo pari a $\gamma(n) \equiv n \cdot g$
- Convesso: un gap di lunghezza n ha costo pari a $\gamma(n)$ tale che

$$\gamma(i + 1) - \gamma(i) \leq \gamma(i) - \gamma(i - 1)$$

- Affine: un gap di lunghezza n ha costo pari a $\gamma(n) \equiv g_o + (n - 1) \cdot g_e$

L'EQUAZIONE DI RICORRENZA

$$A(i, j) = \max \begin{cases} A(i - 1, j - 1) + \text{score}(\sigma_1[i], \sigma_1[j]) \\ \max\{A(k, j) - \gamma(i - k) \mid k = 0, \dots, i - 1\} \\ \max\{A(i, k) - \gamma(j - k) \mid k = 0, \dots, j - 1\} \end{cases}$$

STORIA

- Needleman e Wunsch (1970) allineamento globale
- Smith e Waterman (1981) allineamento locale

MULTIPLE ALIGNMENT

FAINT SIMILARITY STRONG (LOCAL) SIMILARITY

 S_1 *NGPEVRELW* S_2 *ARDIWA* S_3 *QARESIYA* S_4 *VRESLWS* S_5 *WYVRDASLWS*

MULTIPLE ALIGNMENT

FAINT SIMILARITY STRONG (LOCAL) SIMILARITY

S₁	<i>N</i>	<i>G</i>	<i>P</i>	<i>E</i>	<i>V</i>	<i>R</i>	<i>E</i>	–	–	<i>L</i>	<i>W</i>	–
S₂	–	–	–	–	<i>A</i>	<i>R</i>	<i>D</i>	–	–	<i>I</i>	<i>W</i>	<i>A</i>
S₃	–	–	–	<i>Q</i>	<i>A</i>	<i>R</i>	<i>E</i>	–	<i>S</i>	<i>I</i>	<i>Y</i>	<i>A</i>
S₄	–	–	–	–	<i>V</i>	<i>R</i>	<i>E</i>	–	<i>S</i>	<i>L</i>	<i>W</i>	<i>S</i>
S₅	–	–	<i>W</i>	<i>Y</i>	<i>V</i>	<i>R</i>	<i>D</i>	<i>A</i>	<i>S</i>	<i>L</i>	<i>W</i>	<i>S</i>

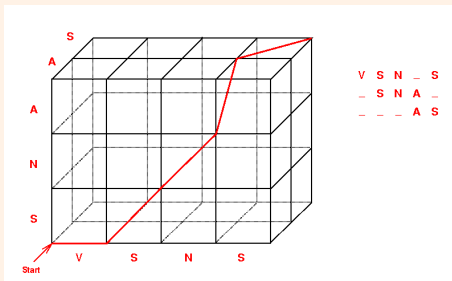
MULTIPLE ALIGNMENT

DEFINITION

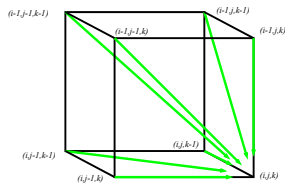
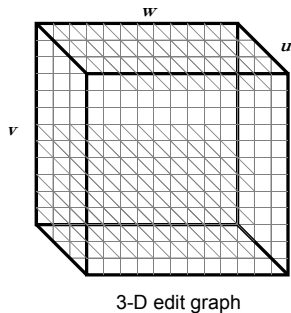
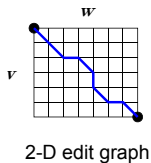
Sum-of-pairs score of a Multiple Sequence Alignment A :

$$Sc(A) = \sum_{i < j} Sc(P_{i,j}(A))$$

DYNAMIC PROGRAMMING?



MULTIPLE ALIGNMENT



MULTIPLE ALIGNMENT

HEURISTICS

Fast algorithms, that do not necessarily find the mathematically optimal alignment.

Two groups $\mathcal{G}_1, \mathcal{G}_2$ of sequences aligned by A_1, A_2 can be aligned like two single sequences if A_1, A_2 remain unchanged.

APPROACH: PROGRESSIVE ALIGNMENT

Align successively sequences and groups of previously aligned sequences, until all sequences are aligned in one MSA.

MULTIPLE ALIGNMENT

PROGRESSIVE ALIGNMENT: EXAMPLE

S₁	<i>E</i>	<i>V</i>	<i>R</i>	<i>E</i>	–	<i>V</i>	<i>W</i>	–
S₂	–	<i>A</i>	<i>R</i>	<i>D</i>	–	<i>T</i>	<i>W</i>	<i>A</i>
S₃	<i>Q</i>	<i>A</i>	<i>R</i>	<i>E</i>	<i>S</i>	<i>I</i>	<i>Y</i>	<i>A</i>
S₄	–	–	<i>R</i>	<i>E</i>	<i>S</i>	<i>L</i>	<i>W</i>	<i>S</i>
S₅	<i>R</i>	<i>E</i>	<i>W</i>	<i>S</i>	<i>L</i>	<i>W</i>	<i>S</i>	
S₆	<i>R</i>	<i>E</i>	<i>Y</i>	<i>S</i>	–	–	<i>S</i>	

MULTIPLE ALIGNMENT

PROGRESSIVE ALIGNMENT: EXAMPLE

S₁	<i>E</i>	<i>V</i>	<i>R</i>	<i>E</i>	–	–	<i>V</i>	<i>W</i>	–
S₂	–	<i>A</i>	<i>R</i>	<i>D</i>	–	–	<i>T</i>	<i>W</i>	<i>A</i>
S₃	<i>Q</i>	<i>A</i>	<i>R</i>	<i>E</i>	–	<i>S</i>	<i>I</i>	<i>Y</i>	<i>A</i>
S₄	–	–	<i>R</i>	<i>E</i>	–	<i>S</i>	<i>L</i>	<i>W</i>	<i>S</i>
S₅	–	–	<i>R</i>	<i>E</i>	<i>W</i>	<i>S</i>	<i>L</i>	<i>W</i>	<i>S</i>
S₆	–	–	<i>R</i>	<i>E</i>	<i>Y</i>	<i>S</i>	–	–	<i>S</i>

MULTIPLE ALIGNMENT

GENERAL STRATEGY

- Construct phylogenetic tree of input sequences S_1, \dots, S_n ('guide tree').
- Traverse T from leaves to root
- At every inner node, construct profile alignments of sequences corresponding to daughter nodes

IMPLEMENTATION

Clustal Ω

DOT-PLOT

	A	A	T	C	T	T	C	A	G	C	G	T	A	T	T	G	C	T
A	1	1	0	0	0	0	0	1	0	0	0	0	1	0	0	0	0	0
T	0	0	1	0	1	1	0	0	0	0	0	1	0	1	1	0	0	1
C	0	0	0	1	0	0	1	0	0	1	0	0	0	0	0	0	1	0
T	0	0	1	0	1	1	0	0	0	0	0	1	0	1	1	0	0	1
T	0	0	1	0	1	1	0	0	0	0	0	1	0	1	1	0	0	1
A	1	1	0	0	0	0	0	1	0	0	0	0	1	0	0	0	0	0
G	0	0	0	0	0	0	0	0	1	0	1	0	0	0	0	1	0	0
C	0	0	0	1	0	0	1	0	0	1	0	0	0	0	0	0	1	0
C	0	0	0	1	0	0	1	0	0	1	0	0	0	0	0	0	1	0
G	0	0	0	0	0	0	0	0	1	0	1	0	0	0	0	1	0	0
G	0	0	0	0	0	0	0	0	1	0	1	0	0	0	0	1	0	0
A	1	1	0	0	0	0	0	1	0	0	0	0	1	0	0	0	0	0
G	0	0	0	0	0	0	0	0	1	0	1	0	0	0	0	1	0	0
T	0	0	0	0	0	0	0	0	1	0	1	0	0	0	0	1	0	0
T	0	0	1	0	1	1	0	0	0	0	0	1	0	1	1	0	0	1
A	1	1	0	0	0	0	0	1	0	0	0	0	1	0	0	0	0	0
T	0	0	1	0	1	1	0	0	0	0	0	1	0	1	1	0	0	1
T	0	0	1	0	1	1	0	0	0	0	0	1	0	1	1	0	0	1

DOT-PLOT

	A	A	T	C	T	T	C	A	G	C	G	T	A	T	T	G	C	T
A	1	1	0	0	0	0	0	1	0	0	0	0	1	0	0	0	0	0
T	0	0	1	0	1	1	0	0	0	0	0	1	0	1	1	0	0	1
C	0	0	0	1	0	0	1	0	0	1	0	0	0	0	0	0	0	1
T	0	0	1	0	1	1	0	0	0	0	0	1	0	1	1	0	0	1
T	0	0	1	0	1	1	0	0	0	0	0	1	0	1	1	0	0	1
A	1	1	0	0	0	0	0	1	0	0	0	0	1	0	0	0	0	0
G	0	0	0	0	0	0	0	0	1	0	1	0	0	0	0	0	1	0
C	0	0	0	1	0	0	1	0	0	1	0	0	0	0	0	0	0	1
C	0	0	0	1	0	0	1	0	0	1	0	0	0	0	0	0	0	1
G	0	0	0	0	0	0	0	0	1	0	1	0	0	0	0	0	1	0
G	0	0	0	0	0	0	0	0	1	0	1	0	0	0	0	0	1	0
A	1	1	0	0	0	0	0	1	0	0	0	0	1	0	0	0	0	0
G	0	0	0	0	0	0	0	0	1	0	1	0	0	0	0	0	1	0
G	0	0	0	0	0	0	0	0	1	0	1	0	0	0	0	0	1	0
T	0	0	1	0	1	1	0	0	0	0	0	1	0	1	1	0	0	1
A	1	1	0	0	0	0	0	1	0	0	0	0	1	0	0	0	0	0
T	0	0	1	0	1	1	0	0	0	0	0	1	0	1	1	0	0	1
T	0	0	1	0	1	1	0	0	0	0	0	1	0	1	1	0	0	1

LE DIAGONALI DI 1 ...

... possono essere molto utili!

BLAST

IDEA

Perché non usare “pezzi” di sequenza uguali *pre-computati*?

BLAST

IDEA

Perché non usare “pezzi” di sequenza uguali *pre-computati*?

DAL CAPITOLO 16 DELL'NCBI HANDBOOK

How BLAST Works: The Basics

The BLAST algorithm is a heuristic program, which means that it relies on some smart shortcuts to perform the search faster. BLAST performs "local" alignments.

...

When a query is submitted via one of the BLAST Web pages, the sequence, plus any other input information such as the database to be searched, word size, expect value, and so on, are fed to the algorithm [http://www.ncbi.nlm.nih.gov/Education/BLASTinfo/BLAST_algorithm.html] on the BLAST server. BLAST works by first making a **look-up table of all the “words”** (short subsequences, which for proteins the default is three letters) and **“neighboring words”**, i.e., similar words in the query sequence. The sequence database is then scanned for these “hot spots”. When a match is identified, it is used to initiate gap-free and gapped extensions of the “word”.

...

BLAST STAGES:

- FIRST STAGE:** Identify exact matches of length W (default $W=3$) between the query and the sequences in the database
- SECOND STAGE:** Extend the match in both directions in an attempt to boost the alignment score (insertions and deletions are not considered)
- THIRD STAGE:** If a high-scoring ungapped alignment is found: Perform a gapped local alignment using dynamic programming

INDEXING

BUILDING AN INDEX FOR A REFERENCE

What is an **Index**?

INDEXING

DEFINITION (WIKIPEDIA)

- A Lookup table, a data structure, usually an array or associative array, often used to replace a runtime computation with a simpler array indexing operation
- Array index, an integer pointer (into an array data structure) that identifies an element of the array
- A key in an associative array
- Database index, a data structure that improves the speed of data retrieval operations on a database table
- Index mapping, mapping of raw data, used directly as in array index, for an array
- Index register, a processor register used for modifying operand addresses during the run of a program
- The dataset maintained by search engine indexing
- ...

INDEXING

BUT ...

... the same idea can be used to build a lookup table of the reference when searching for short strings!

READS AS NUMBERS

DEFINITION

- $\Sigma_{dna} = \{a, c, g, t\}$ alphabet;
- \mathbf{P} pattern and $|\mathbf{P}| = m$;
- \mathbf{T} text and $|\mathbf{T}| = n$;

T_s the m -character string $T[s, \dots, s + m - 1]$.

STRINGS IN Σ_{dna} are NUMBERS IN BASE 4

$$\mathbf{P} \rightsquigarrow p \quad \mathbf{T}_s \rightsquigarrow t_s$$

$$a \equiv 0 \quad c \equiv 1 \quad g \equiv 2 \quad t \equiv 3$$

INDEXING

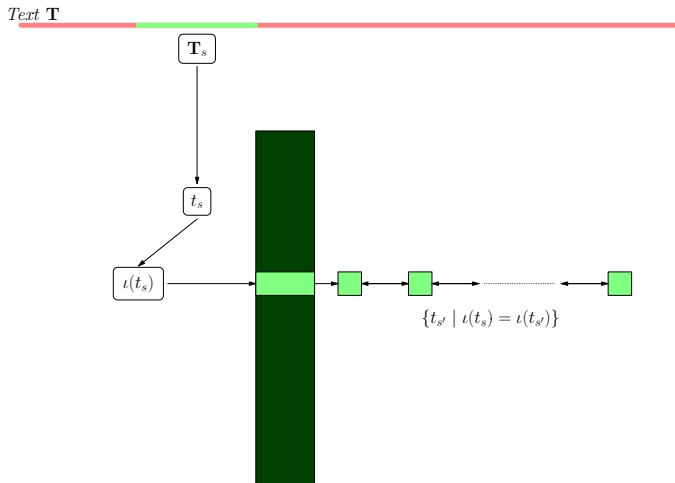
CAREFUL!

4^m is BIG

USE AN *index*

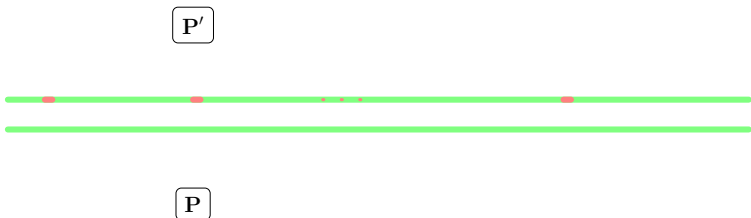
A function ι mapping a—*large*—domain in a—*small*—array

INDEXING



THE PROBLEMS

Mismatches (Insertions/Deletions)



- 1 We must search for a string **and its variants**
- 2 The index should allow to speed-up search. It must be **structured**

UN PASSO *indietro*

DEFINITION (RICERCA ESATTA)

Dato un *pattern* P ed un testo T , trova tutte le occorrenze esatte di P in T .

UN PASSO *indietro*

DEFINITION (RICERCA ESATTA)

Dato un *pattern* P ed un testo T , trova tutte le occorrenze esatte di P in T .

APPLICAZIONI?

Pattern corti.

Pattern lunghi se osservo che sotto-pattern corti *devono* occorrere in modo esatto.

UN PASSO *indietro*

DEFINITION (RICERCA ESATTA)

Dato un *pattern* P ed un testo T , trova tutte le occorrenze esatte di P in T .

APPLICAZIONI?

Pattern corti.

Pattern lunghi se osservo che sotto-pattern corti *devono* occorrere in modo esatto.

COMPLESSITÀ?

UN PASSO *indietro*

DEFINITION (RICERCA ESATTA)

Dato un *pattern* P ed un testo T , trova tutte le occorrenze esatte di P in T .

APPLICAZIONI?

Pattern corti.

Pattern lunghi se osservo che sotto-pattern corti *devono* occorrere in modo esatto.

COMPLESSITÀ?

$$|P| + |T|$$

UN PASSO *indietro*

DEFINITION (RICERCA ESATTA)

Dato un *pattern* P ed un testo T , trova tutte le occorrenze esatte di P in T .

APPLICAZIONI?

Pattern corti.

Pattern lunghi se osservo che sotto-pattern corti *devono* occorrere in modo esatto.

COMPLESSITÀ?

$$|P| + |T|$$

Si può fare di meglio?

DUE OSSERVAZIONI

LEMMA

Se P occorre in T a distanza (Hamming/Levensthein) d , allora una sotto-stringa di P di lunghezza almeno $|P|/(d+1)$ occorre in modo esatto.

DUE OSSERVAZIONI

LEMMA

Se P occorre in T a distanza (Hamming/Levensthein) d , allora una sotto-stringa di P di lunghezza almeno $|P|/(d+1)$ occorre in modo esatto.

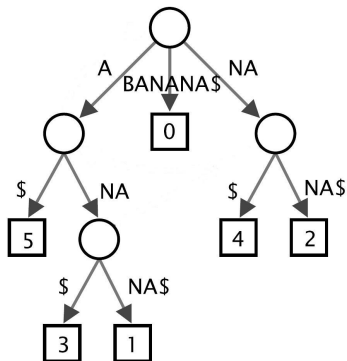
IL TESTO PUÒ ESSERE MEMORIZZATO IN UNA *struttura dati*

- ⇒ non sono obbligato a leggere tutto il testo per ogni pattern
- ⇒ posso cercare P ad un costo $O(|P|)$!

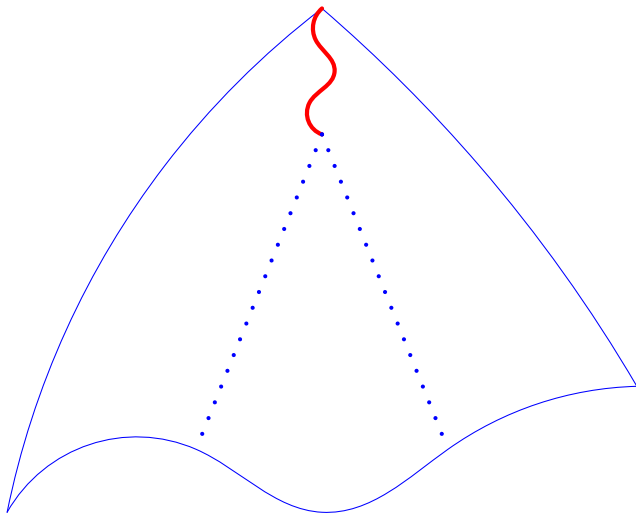
SUFFIX TREES

IDEA

- memorizziamo tutti i suffissi di T
- evitiamo di memorizzare due volte le stesse stringhe (prefissi)



B	A	N	A	N	A	\$
0	1	2	3	4	5	



SUFFIX ARRAYS

i	$A_{SA}[i]$
0	acaaacatat\$
1	caaacatat\$
2	aaacatat\$
3	aacatat\$
4	acatat\$
5	catat\$
6	atat\$
7	tat\$
8	at\$
9	t\$
10	\$

$$T = \text{acaaacatat} \quad P = \text{ata}$$

SUFFIX ARRAYS

i	SA	$A_{SA[i]}$
0	2	aaacatat\$
1	3	aacatat\$
2	0	acaaacatat\$
3	4	acatat\$
4	6	atat\$
5	8	at\$
6	1	caaacatat\$
7	5	catat\$
8	7	tat\$
9	9	t\$
10	10	\$

$$T = \text{acaaacatat} \quad P = \text{ata}$$

SUFFIX ARRAYS

i	SA	$A_{SA[i]}$
0	2	aaacatat\$
1	3	aacatat\$
2	0	acaaacatat\$
3	4	acatat\$
4	6	atat\$
5	8	at\$
6	1	caaacatat\$
7	5	catat\$
8	7	tat\$
9	9	t\$
10	10	\$

$T = \text{acaaacatat\$}$ $P = \text{ata}$

SUFFIX ARRAYS

i	SA	$A_{SA[i]}$
0	2	aaacatat\$
1	3	aacatat\$
2	0	acaaacatat\$
3	4	acatat\$
4	6	atat\$
5	8	at\$
6	1	caaacatat\$
7	5	catat\$
8	7	tat\$
9	9	t\$
10	10	\$

$T = \text{acaaacatat}\$$ $P = \text{ata}$

SUFFIX ARRAYS

i	SA	$A_{SA[i]}$
0	2	aaacatat\$
1	3	aacatat\$
2	0	acaaacatat\$
3	4	acatat\$
4	6	atat\$
5	8	at\$
6	1	caaacatat\$
7	5	catat\$
8	7	tat\$
9	9	t\$
10	10	\$

$T = \text{acaaacatat\$}$ $P = \text{ata}$

SUFFIX ARRAYS

i	SA	$A_{SA[i]}$
0	2	aaacatat\$
1	3	aacatat\$
2	0	acaaacatat\$
3	4	acatat\$
4	6	atat\$
5	8	at\$
6	1	caaacatat\$
7	5	catat\$
8	7	tat\$
9	9	t\$
10	10	\$

$T = \text{acaaacatat}\$$ $P = \text{ata}$

ENHANCED SUFFIX ARRAYS

I SUFFIX ARRAY SONO MIGLIORABILI?

- complessità della costruzione
- complessità della ricerca
- informazioni di supporto (*enhancement*)
- ... compressione: *si possono comprimere contemporaneamente indice e testo?*

LA TRASFORMATATA DI BURROWS-WHEELER

PASSI PER IL CALCOLO

- Considero i primi caratteri
- Ordino lessicograficamente la parte azzurra
- Ottengo una permutazione della stringa di partenza

```
s wiss·miss·missing
w iss·miss·missing
i ss·miss·missing
s s·miss·missing
s ·miss·missing
· miss·missing
m iss·missing
i ss·missing
s s·missing
s ·missing
· missing
m issing
i ssing
s sing
s ing
i ng
n g
g
```


LA TRASFORMATA DI BURROWS-WHEELER

PASSI PER IL CALCOLO

- Considero i primi caratteri
- Ordino lessicograficamente la parte azzurra
- Ottengo una permutazione della stringa di partenza

```
g  
s ·miss·missing  
s ·missing  
n g  
s ing  
w iss·miss·missing  
m iss·missing  
m issing  
· miss·missing  
· missing  
i ng  
s s·miss·missing  
s s·missing  
i ss·miss·missing  
s sing  
i ss·missing  
i ssing  
s wiss·miss·missing
```

LA TRASFORMATATA DI BURROWS-WHEELER

PASSI PER IL CALCOLO

- Considero i primi caratteri
- Ordino lessicograficamente la parte azzurra
- Ottengo una permutazione della stringa di partenza

g
s
s
n
s
w
m
m
.
.
i
s
s
i
s
i
i
s

LA TRASFORMATATA DI BURROWS-WHEELER

- La BWT è reversibile: dalla stringa trasformata è possibile ricostruire il testo originale ...
- ... ma la stringa trasformata è più “facile” da comprimere in quanto “localmente omogenea”.
- Si dimostra che la BWT permette di ottenere una compressione elevata usando compressori “semplici” cioè che non guardano il contesto dei simboli.
- Equivalentemente, comprimendo $BWT(S)$ fino alla sua entropia H_0 e equivalente a comprimere S fino a H_k per ogni $k \geq 0$.