

Avventure Algoritmiche

Bioinformatica e Sequenziamento di Genomi - II

Alberto Policriti



Udine, 30.10.2019

Questioni di base

1. cosa significa *allineare* (due o più sequenze, e.g. di DNA)?
2. perché il problema si presta ad illustrare *algoritmi*?
3. quali obiettivi abbiamo?
4. che dati utilizziamo per fare esperimenti?

Allineamento

Il significato (operativo) del termine allineare

Allineare due stringhe significa ...

... convertire la prima nella seconda

Il significato (operativo) del termine allineare

Allineare due stringhe significa ...

... convertire la prima nella seconda

Example

$\sigma_1 \equiv$ a c g t c a t c a

$\sigma_2 \equiv$ t a a g t g t c a

Che cosa intendiamo per convertire una sequenza in un'altra?

Definizione

Un allineamento di σ_1 e σ_2 è una *stringa* nell'alfabeto

$\{\mathbf{m}, \mathbf{i}, \mathbf{d}, \mathbf{s}\}$

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

Che cosa intendiamo per convertire una sequenza in un'altra?

Definizione

Un allineamento di σ_1 e σ_2 è una *stringa* nell'alfabeto

$$\{\mathbf{m}, \mathbf{i}, \mathbf{d}, \mathbf{s}\}$$

Example

$$\begin{array}{rcccccccccc} \sigma_1 \equiv & & a & c & g & t & c & a & t & c & a \\ & & \mathbf{i} & \mathbf{m} & \mathbf{s} & \mathbf{m} & \mathbf{m} & \mathbf{s} & \mathbf{d} & \mathbf{m} & \mathbf{m} & \mathbf{m} \\ \sigma_2 \equiv & t & a & a & g & t & g & & t & c & a \end{array}$$

Un allineamento è una stringa
allineamento \Leftrightarrow edit-string

In effetti è un (semplice, lineare) “*programma*” che converte σ_1 in σ_2 (o viceversa).

Che cosa intendiamo per convertire una sequenza in un'altra?

Definizione

Un allineamento di σ_1 e σ_2 è una *stringa* nell'alfabeto

$$\{\mathbf{m}, \mathbf{i}, \mathbf{d}, \mathbf{s}\}$$

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

Due problemi

1. *trovare* un allineamento
2. stabilire se quello trovato è *ottimo* (in che senso?)

Ordiniamo Allineamenti: Costo e Punteggio

Possiamo associare un **costo** ad una edit-string (vogliamo ↓)

Possiamo associare un **punteggio** ad una edit-string (vogliamo ↑)

Ordiniamo Allineamenti: Costo e Punteggio

Possiamo associare un **costo** ad una edit-string (vogliamo ↓)

Possiamo associare un **punteggio** ad una edit-string (vogliamo ↑)

Assegnamo un Punteggio

match (buono): punteggio *alto* (e.g. 1)

substitution (cattivo): punteggio basso (e.g. 0)

insertion/deletion (cattivo): punteggio basso (e.g. 0)

Ordiniamo Allineamenti: Costo e Punteggio

Possiamo associare un **costo** ad una edit-string (vogliamo ↓)

Possiamo associare un **punteggio** ad una edit-string (vogliamo ↑)

Example

Purine (a,g) e Pirimidine (c,t) sono simili, quindi una possibile (non-banale) **matrice di punteggi** è la seguente:

	<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

Ordiniamo Allineamenti: Costo e Punteggio

Possiamo associare un **costo** ad una edit-string (vogliamo ↓)

Possiamo associare un **punteggio** ad una edit-string (vogliamo ↑)

Example

Penalizzazione dei gap (Affine). In effetti ci piacerebbe anche:

penalizzare l'**apertura** di un gap

e

penalizzare una sua **estensione**

... dobbiamo integrare la funzione di costo nel meccanismo ricorsivo.

Globale/Locale

La edit-string converte l'*intera* σ_2 nell'*intera* σ_1 .

La edit-string converte l'*intera* σ_2 in *una porzione* di σ_1 .

Distanza

Pensiamo al **costo** come una *distanza* tra stringhe.

Example

Distanza di Levenstein: costo 0 per **m**, 1 per gli altri.

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

Distanza di Hamming: **i** e **d** non sono ammesse.

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

La distanza di Hamming può essere usata esclusivamente su stringhe di uguale lunghezza.

Che distanza devo usare?

Stringhe lunghe \Leftrightarrow Levensthein

Stringhe corte \Leftrightarrow Hamming

Algoritmi di Allineamento

Idea!

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

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: $\eta[k + 1] \in \{\mathbf{m}, \mathbf{s}, \mathbf{i}, \mathbf{d}\}$.

Idea!

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

- 1,2** Se $\eta[k + 1] \in \{\mathbf{m}, \mathbf{s}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i]$ e $\sigma_2[1, \dots, j]$.
- 3** Se $\eta[k + 1] \in \{\mathbf{i}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i - 1]$ e $\sigma_2[1, \dots, j]$.
- 4** Se $\eta[k + 1] \in \{\mathbf{d}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i]$ e $\sigma_2[1, \dots, j - 1]$.

Idea!

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

- 1,2** Se $\eta[k + 1] \in \{\mathbf{m}, \mathbf{s}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i]$ e $\sigma_2[1, \dots, j]$.
- 3** Se $\eta[k + 1] \in \{\mathbf{i}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i - 1]$ e $\sigma_2[1, \dots, j]$.
- 4** Se $\eta[k + 1] \in \{\mathbf{d}\}$, allora $\eta[1, \dots, k + 1]$ allinea $\sigma_1[1, \dots, i]$ e $\sigma_2[1, \dots, j - 1]$.

Determino $\eta[k + 1]$ scegliendo quello che mi da il costo minimo.

Algoritmo ricorsivo

... esponenziale

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 $\min(c_m, c_s, c_i, c_d)$;

$$T_{\text{allineamento_exp}}(1) = O(1)$$

$$T_{\text{allineamento_exp}}(n) = 3T_{\text{allineamento_exp}}(n - 1)$$

$$T_{all_exp}(1) = O(1)$$

$$T_{all_exp}(i + 0) = O(1)$$

$$T_{all_exp}(0 + j) = O(1)$$

$$T_{all_exp}(i + j) = 2T_{all_exp}(i + j - 1) + T_{all_exp}(i + j - 2)$$

... ricordiamo che lavoriamo nel caso *pessimo*

Vedremo dopo. Nel frattempo osserviamo che c'è un *trade-off* tempo spazio.

Ci sono al più $|\sigma_1| \cdot |\sigma_2|$ valori per le chiamate ricorsive ...
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 $\sigma_2[1, \dots, j]$

Ci sono al più $|\sigma_1| \cdot |\sigma_2|$ valori per le chiamate ricorsive ...

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	j	
ϵ			
		\vdots	
i	...	$d_L(i, j)$...
		\vdots	

Algorithm 2: *allineamento_prog_din*(σ_1, σ_2, i, j)

```
for  $i = 0; i \leq |\sigma_1|; i++$  do // inizializzo (la prima colonna di  $A$ )  
   $A(i, 0) = i$ 
```

```
for  $j = 0; j \leq |\sigma_2|; j++$  do // inizializzo (la prima riga di  $A$ )  
   $A(0, j) = j$ 
```

```
// completo  $A$  procedendo  $\downarrow$  sulle colonne)
```

```
for  $j = 0; j \leq |\sigma_2|; j++$  do  
  for  $i = 0; i \leq |\sigma_1|; i++$  do  
    if  $\sigma_1[i] == \sigma_2[j]$  then  
       $A(i, j) = \min\{A(i-1, j-1), A(i-1, j) + 1, A(i, j-1) + 1\};$   
    else  
       $A(i, j) = \min\{A(i-1, j-1) + 1, A(i-1, j) + 1, A(i, j-1) + 1\};$ 
```

Example

	ϵ	a	c	g	t	c	a	t	c	a
ϵ	0	1	2	3	4	5	6	7	8	9
t	1	1	2	3	3	4	5	6	7	8
a	2	1	2	3	4	4	4	5	6	7
a	3	2	2	3	4	5	4	5	6	6
g	4	3	3	2	4	5	5	5	6	7
t	5	4	4	3	2	3	4	5	6	7
g	6	5	5	4	3	3	4	5	6	7
t	7	6	6	5	4	4	4	4	5	6
c	8	7	6	6	5	4	5	5	4	5
a	9	8	7	7	6	5	4	5	5	4

Example

	ϵ	a	c	g	t	c	a	t	c	a
ϵ	0	1	2	3	4	5	6	7	8	9
t	1	1	2	3	3	4	5	6	7	8
a	2	1	2	3	4	4	4	5	6	7
a	3	2	2	3	4	5	4	5	6	6
g	4	3	3	2	4	5	5	5	6	7
t	5	4	4	3	2	3	4	5	6	7
g	6	5	5	4	3	3	4	5	6	7
t	7	6	6	5	4	4	4	4	5	6
c	8	7	6	6	5	4	5	5	4	5
a	9	8	7	7	6	5	4	5	5	4

$\sigma_1 \equiv$ *a c g t c a t c a*

i m s m m s d m m m $\Rightarrow d_L(\sigma_1, \sigma_2) = 4.$

$\sigma_2 \equiv$ *t a a g t g t c a*

Example

	ϵ	a	c	g	t	c	a	t	c	a
ϵ	0	1	2	3	4	5	6	7	8	9
t	1	1	2	3	3	4	5	6	7	8
a	2	1	2	3	4	4	4	5	6	7
a	3	2	2	3	4	5	4	5	6	6
g	4	3	3	2	4	5	5	5	6	7
t	5	4	4	3	2	3	4	5	6	7
g	6	5	5	4	3	3	4	5	6	7
t	7	6	6	5	4	4	4	4	5	6
c	8	7	6	6	5	4	5	5	4	5
a	9	8	7	7	6	5	4	5	5	4

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

Example

	ϵ	a	c	g	t	c	a	t	c	a
ϵ	0	1	2	3	4	5	6	7	8	9
t	1	1	2	3	3	4	5	6	7	8
a	2	1	2	3	4	4	4	5	6	7
a	3	2	2	3	4	5	4	5	6	6
g	4	3	3	2	4	5	5	5	6	7
t	5	4	4	3	2	3	4	5	6	7
g	6	5	5	4	3	3	4	5	6	7
t	7	6	6	5	4	4	4	4	5	6
c	8	7	6	6	5	4	5	5	4	5
a	9	8	7	7	6	5	4	5	5	4

$\sigma_1 \equiv$ **a** **c** **g** **t** **c** **a** **t** **c** **a**

i **m** **s** **m** **m** **s** **d** **m** **m** **m** $\Rightarrow d_L(\sigma_1, \sigma_2) = 4.$

$\sigma_2 \equiv$ **t** **a** **a** **g** **t** **g** **t** **c** **a**

Example

	ϵ	a	c	g	t	c	a	t	c	a
ϵ	0	1	2	3	4	5	6	7	8	9
t	1	1	2	3	3	4	5	6	7	8
a	2	1	2	3	4	4	4	5	6	7
a	3	2	2	3	4	5	4	5	6	6
g	4	3	3	2	4	5	5	5	6	7
t	5	4	4	3	2	3	4	5	6	7
g	6	5	5	4	3	3	4	5	6	7
t	7	6	6	5	4	4	4	4	5	6
c	8	7	6	6	5	4	5	5	4	5
a	9	8	7	7	6	5	4	5	5	4

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

Example

	ε	a	c	g	t	c	a	t	c	a
ε	0	1	2	3	4	5	6	7	8	9
t	1	1	2	3	3	4	5	6	7	8
a	2	1	2	3	4	4	4	5	6	7
a	3	2	2	3	4	5	4	5	6	6
g	4	3	3	2	4	5	5	5	6	7
t	5	4	4	3	2	3	4	5	6	7
g	6	5	5	4	3	3	4	5	6	7
t	7	6	6	5	4	4	4	4	5	6
c	8	7	6	6	5	4	5	5	4	5
a	9	8	7	7	6	5	4	5	5	4

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

Example

	ε	a	c	g	t	c	a	t	c	a
ε	0	1	2	3	4	5	6	7	8	9
t	1	1	2	3	3	4	5	6	7	8
a	2	1	2	3	4	4	4	5	6	7
a	3	2	2	3	4	5	4	5	6	6
g	4	3	3	2	4	5	5	5	6	7
t	5	4	4	3	2	3	4	5	6	7
g	6	5	5	4	3	3	4	5	6	7
t	7	6	6	5	4	4	4	4	5	6
c	8	7	6	6	5	4	5	5	4	5
a	9	8	7	7	6	5	4	5	5	4

$\sigma_1 \equiv$ a c g t c a t c a

 i m s m m s d m m m $\Rightarrow d_L(\sigma_1, \sigma_2) = 4.$

$\sigma_2 \equiv$ t a a g t g t c a

Example

	ϵ	a	c	g	t	c	a	t	c	a
ϵ	0	1	2	3	4	5	6	7	8	9
t	1	1	2	3	3	4	5	6	7	8
a	2	1	2	3	4	4	4	5	6	7
a	3	2	2	3	4	5	4	5	6	6
g	4	3	3	2	4	5	5	5	6	7
t	5	4	4	3	2	3	4	5	6	7
g	6	5	5	4	3	3	4	5	6	7
t	7	6	6	5	4	4	4	4	5	6
c	8	7	6	6	5	4	5	5	4	5
a	9	8	7	7	6	5	4	5	5	4

$\sigma_1 \equiv$ a c g t c a t c a

 i m s m m s d m m m $\Rightarrow d_L(\sigma_1, \sigma_2) = 4.$

$\sigma_2 \equiv$ t a a g t g t c a

Example

	ε	a	c	g	t	c	a	t	c	a
ε	0	1	2	3	4	5	6	7	8	9
t	1	1	2	3	3	4	5	6	7	8
a	2	1	2	3	4	4	4	5	6	7
a	3	2	2	3	4	5	4	5	6	6
g	4	3	3	2	4	5	5	5	6	7
t	5	4	4	3	2	3	4	5	6	7
g	6	5	5	4	3	3	4	5	6	7
t	7	6	6	5	4	4	4	4	5	6
c	8	7	6	6	5	4	5	5	4	5
a	9	8	7	7	6	5	4	5	5	4

$\sigma_1 \equiv$ *a c g t c a t c a*

i m s m m s d m m m $\Rightarrow d_L(\sigma_1, \sigma_2) = 4.$

$\sigma_2 \equiv$ *t a a g t g t c a*

Example

	ε	a	c	g	t	c	a	t	c	a
ε	0	1	2	3	4	5	6	7	8	9
t	1	1	2	3	3	4	5	6	7	8
a	2	1	2	3	4	4	4	5	6	7
a	3	2	2	3	4	5	4	5	6	6
g	4	3	3	2	4	5	5	5	6	7
t	5	4	4	3	2	3	4	5	6	7
g	6	5	5	4	3	3	4	5	6	7
t	7	6	6	5	4	4	4	4	5	6
c	8	7	6	6	5	4	5	5	4	5
a	9	8	7	7	6	5	4	5	5	4

$\sigma_1 \equiv$ *a c g t c a t c a*

i m s m m s d m m m $\Rightarrow d_L(\sigma_1, \sigma_2) = 4.$

$\sigma_2 \equiv$ *t a a g t g t c a*

Example

	ϵ	a	c	g	t	c	a	t	c	a
ϵ	0	1	2	3	4	5	6	7	8	9
t	1	1	2	3	3	4	5	6	7	8
a	2	1	2	3	4	4	4	5	6	7
a	3	2	2	3	4	5	4	5	6	6
g	4	3	3	2	4	5	5	5	6	7
t	5	4	4	3	2	3	4	5	6	7
g	6	5	5	4	3	3	4	5	6	7
t	7	6	6	5	4	4	4	4	5	6
c	8	7	6	6	5	4	5	5	4	5
a	9	8	7	7	6	5	4	5	5	4

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

Example

	ε	a	c	g	t	c	a	t	c	a
ε	0	1	2	3	4	5	6	7	8	9
t	1	1	2	3	3	4	5	6	7	8
a	2	1	2	3	4	4	4	5	6	7
a	3	2	2	3	4	5	4	5	6	6
g	4	3	3	2	4	5	5	5	6	7
t	5	4	4	3	2	3	4	5	6	7
g	6	5	5	4	3	3	4	5	6	7
t	7	6	6	5	4	4	4	4	5	6
c	8	7	6	6	5	4	5	5	4	5
a	9	8	7	7	6	5	4	5	5	4

$\sigma_1 \equiv$ *a c g t c a t c a*

i m s m m s d m m m $\Rightarrow d_L(\sigma_1, \sigma_2) = 4.$

$\sigma_2 \equiv$ *t a a g t g t c a*

- Si può leggere un allineamento sulla matrice da sx a dx e viceversa
- Si può risparmiare spazio (interessante!)
- Diverse regioni della matrice possono apparire inutili ma poi rivelarsi importanti

Domanda:

Cosa ho *veramente* computato?

Risposta:

L'allineamento di un prefisso di σ_1 con un prefisso di σ_2

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

Definizione

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 $\leq d$).

Basta ...

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

Example

	€	a	c	g	t	a	c	g	c	g	t	a	c	a	g	t	a	...
€	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	...
c	1	1	0	1	1	1	0	1	0	1	1	1	0	1	1	1	1	...
g	2	2	1	0	1	2	1	0	1	0	1	2	1	1	1	2	2	...
t	3	3	2	2	0	1	2	1	1	1	0	1	2	2	2	1	2	...
a	4	3	3	3	1	0	1	2	2	2	1	0	1	2	3	2	1	...

Example

	€	a	c	g	t	a	c	g	c	g	t	a	c	a	g	t	a	...
€	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	...
c	1	1	0	1	1	1	0	1	0	1	1	1	0	1	1	1	1	...
g	2	2	1	0	1	2	1	0	1	0	1	2	1	1	1	2	2	...
t	3	3	2	2	0	1	2	1	1	1	0	1	2	2	2	1	2	...
a	4	3	3	3	1	0	1	2	2	2	1	0	1	2	3	2	1	...

Abbiamo introdotto un gap iniziale a costo zero

La soluzione al problema dell'allineamento locale si ottiene semplicemente ragionando sul metodo impiegato per risolvere l'allineamento globale.

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 $\gamma(n) \equiv g_o + (n-1) \cdot g_e$

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

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

Possiamo guardare indietro (nella ricorsione) quanto ci pare!

Obiettivi

Due categorie di obiettivi

- usare dati *interessanti*
- ragionare sulla ricorsione

Caratteristiche

- Il procedimento ricorsivo è inusuale
- Gli studenti tendono a “seguire” le chiamate ricorsive fino alla base (\Rightarrow tabulazione)
- Una volta acquisito è uno strumento potente (torre di Hanoi)
- Consente di analizzare la complessità computazionale mediante **equazioni**:

$$T_{\text{allineamento_exp}}(1) = O(1)$$

$$T_{\text{allineamento_exp}}(n) = 3T_{\text{allineamento_exp}}(n - 1)$$

Tabulazione dei risultati delle chiamate ricorsive

- Non è sempre necessaria
- Non è sempre sensato farla

La programmazione dinamica è un caso particolare: ricorsione “strutturata”

- ci sono *tanti* allineamenti (tanti percorsi per trovarli) ma ...
- ... alla fine ne vogliamo uno (il migliore)!

- ci sono *tanti* allineamenti (tanti percorsi per trovarli) ma ...
- ... alla fine ne vogliamo uno (il migliore)!

L'argomento si presta a:

- ragionare sull'importanza di definire chiaramente il problema:
come definisco il miglior allineamento?
- ragionare come integrare conoscenza di un campo (Biologia)
in uno strumento algoritmico (Computazione): come garantisco che il
miglior allineamento valuti correttamente il costo dei gap?
- considerare la ricorsione come strumento per analizzare lo
spazio delle soluzioni: non mi interessa come ho ottenuto la soluzione ad un
sottoproblema, mi interessa solo che quel sottoproblema ha soluzione.

Analizzare la ricorsione

Tabulare i risultati e scriverli in una matrice per stimarne il numero.

Example

	ε	a	c	g	t	c	a	t	c	a
ε	0	1	2	3	4	5	6	7	8	9
t	1	1	2	3	3	4	5	6	7	8
a	2	1	2	3	4	4	4	5	6	7
a	3	2	2	3	4	5	4	5	6	6
g	4	3	3	2	4	5	5	5	6	7
t	5	4	4	3	2	3	4	5	6	7
g	6	5	5	4	3	3	4	5	6	7
t	7	6	6	5	4	4	4	4	5	6
c	8	7	6	6	5	4	5	5	4	5
a	9	8	7	7	6	5	4	5	5	4

$\sigma_1 \equiv$ *a c g t c a t c a*

i m s m m s d m m m $\Rightarrow d_L(\sigma_1, \sigma_2) = 4.$

$\sigma_2 \equiv$ *t a a g t g* *t c a*