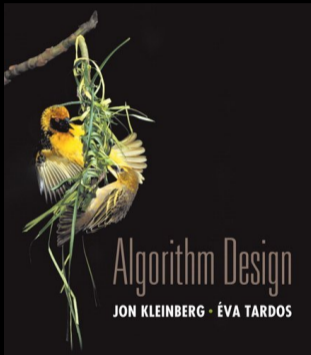


Chapter 6

Dynamic Programming



PEARSON
Addison
Wesley

Paradigmi algoritmici

Greedy. Costruisce la soluzione incrementalmente, ottimizzando qualche criterio locale.

Divide-and-conquer. Divide un problema in diversi sottoproblemi, risolve ogni sottoproblema indipendentemente, e combina le soluzioni ai sottoproblemi per costruire la soluzione al problema originale.

Programmazione dinamica. Divide il problema in una serie di sottoproblemi che si sovrappongono, e costruisce le soluzioni dai sottoproblemi più piccoli a quelli sempre più grandi.

Programmazione Dinamica: storia

Richard Bellman. E' stato il pioniere dello studio sistematico della programmazione dinamica negli anni '50.

Etimologia.

- Programmazione dinamica = pianificazione nel tempo.
programmazione matematica e non programmazione informatica
- Assistant Secretary of Air Force era ostile alla ricerca matematica.
- Bellman pensò ad un nome altisonante per evitarne il confronto.
 - "it's impossible to use dynamic in a pejorative sense"
 - "something not even a Congressman could object to"



Programmazione Dinamica: Applicazioni

Aree.

- Bioinformatica.
- Teoria dei Controlli.
- Teoria dell' Informazione.
- Ricerca Operativa.
- Informatica: teoria, grafica, Intelligenza Artificiale,

Alcni algoritmi di programmazione dinamica famosi.

- Viterbi for hidden Markov models.
- Unix diff for comparing two files.
- Smith-Waterman for sequence alignment.
- Bellman-Ford for shortest path routing in networks.
- Cocke-Kasami-Younger for parsing context free grammars.

Programmazione Dinamica

Algoritmi che vedremo:

- ❑ Weighted Interval Scheduling
- ❑ Segmented Least Squares
- ❑ Knapsack Problem
- ❑ RNA Secondary Structure
- ❑ Sequence Alignment
- ❑ Shortest Paths in a Graph

Esercizio:

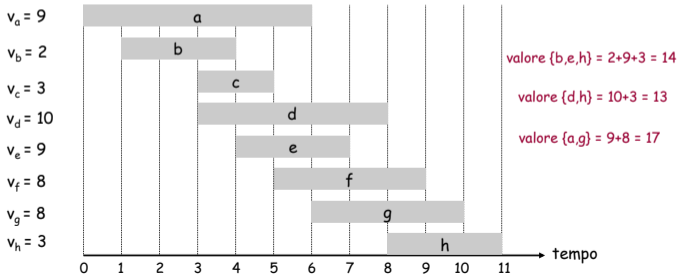
- ❑ Longest Common Subsequence

6.1 Weighted Interval Scheduling

Schedulazione degli intervalli pesati

Problema della schedulazione degli intervalli pesati.

- Job j inizia ad s_j , finisce a f_j , ed ha peso / valore v_j .
- Due job **compatibili** se non hanno intersezione.
- Obiettivo: trovare sottoinsieme con massimo peso di job mutuamente compatibili.



Schedulazione intervalli senza pesi

Ricordiamo. L'algoritmo greedy risolve il problema se tutti i pesi sono 1.

- Considerare job in ordine crescente di tempo di fine.
- Aggiungere job al sottoinsieme della soluzione se è compatibile con i job scelti precedentemente.

Osservazione. L'algoritmo greedy può dare soluzioni non ottime se ci sono pesi arbitrari.

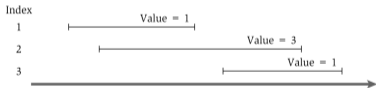
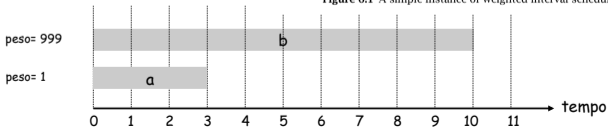


Figure 6.1 A simple instance of weighted interval scheduling.

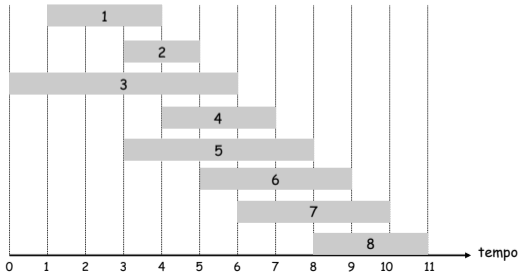


Schedulazione degli intervalli pesati

Notazione. Ordinare job per tempo di fine: $f_1 \leq f_2 \leq \dots \leq f_n$.

Definizione. $p(j)$ = il più grande $i < j$ tale che job i è compatibile con j .

Esempio: $p(8) = 5$, $p(7) = 3$, $p(2) = 0$.



Schedulazione degli intervalli pesati

Notazione. Ordinare job per tempo di fine: $f_1 \leq f_2 \leq \dots \leq f_n$.

Definizione. $p(j)$ = il più grande $i < j$ tale che job i è compatibile con j .

Esempio:

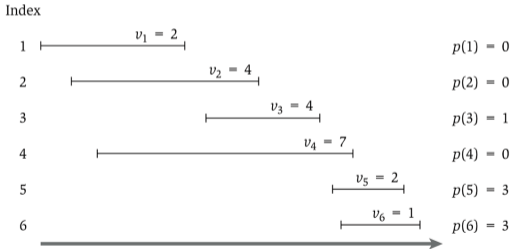


Figure 6.2 An instance of weighted interval scheduling with the functions $p(j)$ defined for each interval j .

Programmazione dinamica: Scelta binaria

Notazione. $OPT(j)$ = valore della soluzione ottimale al problema con i job $1, 2, \dots, j$.

- **Caso 1:** OPT contiene job j .
 - non può contenere job incompatibili $\{ p(j) + 1, p(j) + 2, \dots, j - 1 \}$
 - deve includere soluzione ottimale al problema che consiste dei rimanenti job compatibili $1, 2, \dots, p(j)$
- **Caso 2:** OPT non contiene job j .
 - Deve includere soluzione ottimale al problema che consiste dei rimanenti job compatibili $1, 2, \dots, j-1$

sottostruttura ottima

Programmazione dinamica: Scelta binaria

Notazione. $OPT(j)$ = valore della soluzione ottimale al problema con i job $1, 2, \dots, j$.

- **Caso 1:** OPT contiene job j .
 - non può contenere job incompatibili $\{ p(j) + 1, p(j) + 2, \dots, j - 1 \}$
 - deve includere soluzione ottimale al problema che consiste dei rimanenti job compatibili $1, 2, \dots, p(j)$
- **Caso 2:** OPT non contiene job j .
 - Deve includere soluzione ottimale al problema che consiste dei rimanenti job compatibili $1, 2, \dots, j-1$

sottostruttura ottima

$$OPT(j) = \begin{cases} 0 & \text{se } j=0 \\ \max \{ v_j + OPT(p(j)), OPT(j-1) \} & \text{altrimenti} \end{cases}$$

Schedulazione degli intervalli pesati: Algoritmo di forza bruta

Algoritmo di forza bruta.

Input: $n, s_1, \dots, s_n, f_1, \dots, f_n, v_1, \dots, v_n$

Ordinare job per tempo di fine $f_1 \leq f_2 \leq \dots \leq f_n$.

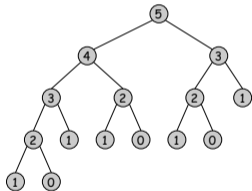
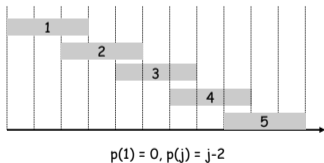
Computare $p(1), p(2), \dots, p(n)$

```
Compute-Opt(j) {  
    if (j = 0)  
        return 0  
    else  
        return max( $v_j + \text{Compute-Opt}(p(j))$ ,  $\text{Compute-Opt}(j-1)$ )  
}
```


Schedulazione degli intervalli pesati: Algoritmo di forza bruta

Osservazione. Algoritmo ricorsivo non è efficiente perchè ci sono molti sottoproblemi ridondanti \Rightarrow complessità esponenziale.

Esempio. Numero di chiamate ricorsive per una famiglia di istanze cresce in modo simile ad i numeri di Fibonacci.



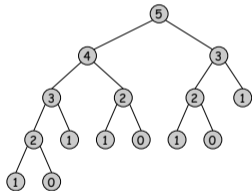
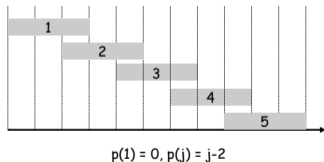
$$T(n) = T(n-1) + T(n-2) + O(1)$$

Schedulazione degli intervalli pesati: Algoritmo di forza bruta

Osservazione. Algoritmo ricorsivo non è efficiente perchè ci sono molti sottoproblemi ridondanti \Rightarrow complessità esponenziale.

Esempio. Numero di chiamate ricorsive per una famiglia di istanze cresce in modo simile ad i numeri di Fibonacci.

i	5	4	3	2	1
	14	8	4	2	0



$$T(n) = T(n-1) + T(n-2) + O(1)$$

$$F_i = F_{i-1} + F_{i-2}$$
$$F_1 = 1 \quad F_0 = 0$$

$$F_i = \frac{\varphi^i - \hat{\varphi}^i}{\sqrt{5}}$$

$$\varphi = \frac{1 + \sqrt{5}}{2} = 1,61803\dots$$

$$\hat{\varphi} = \frac{1 - \sqrt{5}}{2} = -0,61803\dots$$

$$T(n) = \Theta(\varphi^n)$$

Schedulazione degli intervalli pesati: Annotazione

Annotazione (Memoization). Memorizzare i risultati per ogni sottoproblema e verificare quando necessario se un sottoproblema è già stato risolto.

Input: $n, s_1, \dots, s_n, f_1, \dots, f_n, v_1, \dots, v_n$

Ordinare job per tempo di fine $f_1 \leq f_2 \leq \dots \leq f_n$.

Compute $p(1), p(2), \dots, p(n)$

for $j = 1$ to n

$M[j] = \text{empty}$ ← array globale

M-Compute-Opt(j)

 If $j = 0$ then

 Return 0

 Else if $M[j]$ is not empty then

 Return $M[j]$

 Else

 Define $M[j] = \max(v_j + M\text{-Compute-Opt}(p(j)), M\text{-Compute-Opt}(j - 1))$

 Return $M[j]$

 Endif

Schedulazione degli intervalli pesati: Complessità

Claim. L' algoritmo con l'annotazione ha complessità $O(n \log n)$.

- Ordinamento per tempo di fine: $O(n \log n)$.
- Computazione di $p(\cdot)$: $O(n)$ dopo un ordinamento per tempo di inizio.

← esercizio da fare a casa

Schedulazione degli intervalli pesati: Complessità

Claim. L' algoritmo con l' annotazione ha complessità $O(n \log n)$.

- Ordinamento per tempo di fine: $O(n \log n)$.
- Computazione di $p(\cdot)$: $O(n)$ dopo un ordinamento per tempo di inizio.
- $M\text{-Compute-Opt}(j)$: ogni chiamata richiede tempo $O(1)$ e fa una delle due
 - ritorna un valore già calcolato $M[j]$
 - calcola un nuovo valore $M[j]$ e fa due chiamate ricorsive
- Misura del progresso Φ = numero di valori nonvuoti di $M[\cdot]$.
 - inizialmente $\Phi = 0$, poi $\Phi \leq n$.
 - incremento Φ di 1 \Rightarrow al massimo totale di $O(n)$ chiamate ricorsive.
- Tempo totale di $M\text{-Compute-Opt}(n)$ è $O(n)$. ▪

← esercizio da fare a casa

Osservazione. Tempo $O(n)$ se job sono ordinati per tempo di inizio e per tempo di fine.

Annotazione automatica

Annotazione automatica. Molti linguaggi di programmazione funzionali (e.g., Lisp) hanno un supporto built-in support per l'annotazione.

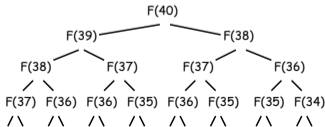
Non succede nei linguaggi imperativi (e.g., Java)?

```
(defun F (n)
  (if
    (<= n 1)
    n
    (+ (F (- n 1)) (F (- n 2)))))
```

Lisp (efficiente)

```
static int F(int n) {
  if (n <= 1) return n;
  else return F(n-1) + F(n-2);
}
```

Java (esponenziale)



Schedulazione degli intervalli pesati: Trovare una soluzione

Domanda. Algoritmo di programmazione dinamica computa il valore ottimale. E se voglio una soluzione?

Risposta. Necessarie operazioni di post-processing.

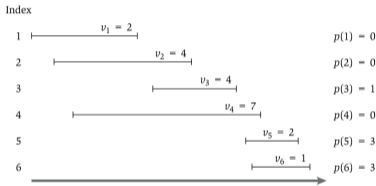


Figure 6.2 An instance of weighted interval scheduling with the functions $p(j)$ defined for each interval j .

$$M = \begin{array}{|c|c|c|c|c|c|c|} \hline & 0 & 1 & 2 & 3 & 4 & 5 & 6 \\ \hline 0 & 0 & 2 & 4 & 6 & 7 & 8 & 8 \\ \hline \end{array}$$

Schedulazione degli intervalli pesati: Trovare una soluzione

Domanda. Algoritmo di programmazione dinamica computa il valore ottimale. E se voglio una soluzione?

Risposta. Necessarie operazioni di post-processing.

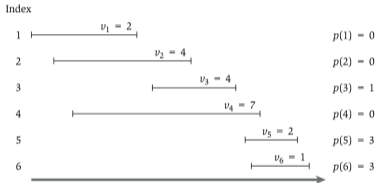


Figure 6.2 An instance of weighted interval scheduling with the functions $p(j)$ defined for each interval j .

E' 6 nella soluzione ottimale?

	0	1	2	3	4	5	6
$M =$	0	2	4	6	7	8	8

Schedulazione degli intervalli pesati: Trovare una soluzione

Domanda. Algoritmo di programmazione dinamica computa il valore ottimale. E se voglio una soluzione?

Risposta. Necessarie operazioni di post-processing.

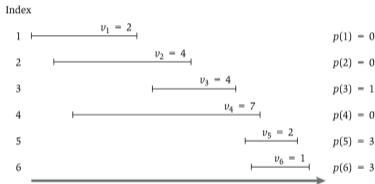


Figure 6.2 An instance of weighted interval scheduling with the functions $p(j)$ defined for each interval j .

E' 6 nella soluzione ottimale? NO

$$OPT(j) = \max \{ v_j + OPT(p(j)), OPT(j-1) \}$$

$$OPT(6) = \max \{ v_6 + OPT(p(6)), OPT(5) \} = \max \{ 1 + OPT(3), OPT(5) \} = \max \{ 1 + 6, 8 \}$$

	0	1	2	3	4	5	6
$M =$	0	2	4	6	7	8	8

Schedulazione degli intervalli pesati: Trovare una soluzione

Domanda. Algoritmo di programmazione dinamica computa il valore ottimale. E se voglio una soluzione?

Risposta. Necessarie operazioni di post-processing.

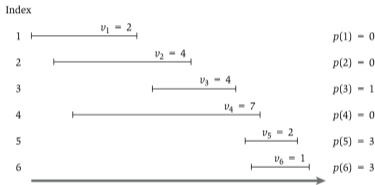


Figure 6.2 An instance of weighted interval scheduling with the functions $p(j)$ defined for each interval j .

E' 5 nella soluzione ottimale?

	0	1	2	3	4	5	6
$M =$	0	2	4	6	7	8	8

Schedulazione degli intervalli pesati: Trovare una soluzione

Domanda. Algoritmo di programmazione dinamica computa il valore ottimale. E se voglio una soluzione?

Risposta. Necessarie operazioni di post-processing.

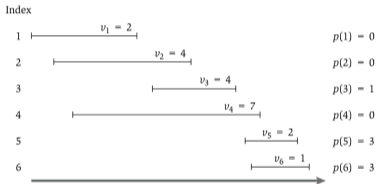


Figure 6.2 An instance of weighted interval scheduling with the functions $p(j)$ defined for each interval j .

E' 5 nella soluzione ottimale? SI

$$M = \begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 1 & 2 & 3 & 4 & 5 & 6 \\ \hline 0 & 2 & 4 & 6 & 7 & 8 & 8 \\ \hline \end{array}$$

$$OPT(j) = \max \{ v_j + OPT(p(j)), OPT(j-1) \}$$

$$OPT(5) = \max \{ v_5 + OPT(p(5)), OPT(4) \} = \max \{ 2 + OPT(3), OPT(4) \} = \max \{ 2 + 6, 7 \}$$

Schedulazione degli intervalli pesati: Trovare una soluzione

Domanda. Algoritmo di programmazione dinamica computa il valore ottimale. E se voglio una soluzione?

Risposta. Necessarie operazioni di post-processing.

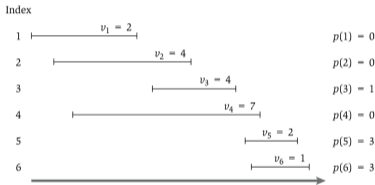


Figure 6.2 An instance of weighted interval scheduling with the functions $p(j)$ defined for each interval j .

E' 3 nella soluzione ottimale?

$$M = \begin{array}{|c|c|c|c|c|c|c|} \hline & 0 & 1 & 2 & 3 & 4 & 5 & 6 \\ \hline & 0 & 2 & 4 & 6 & 7 & 8 & 8 \\ \hline \end{array}$$

Schedulazione degli intervalli pesati: Trovare una soluzione

Domanda. Algoritmo di programmazione dinamica computa il valore ottimale. E se voglio una soluzione?

Risposta. Necessarie operazioni di post-processing.

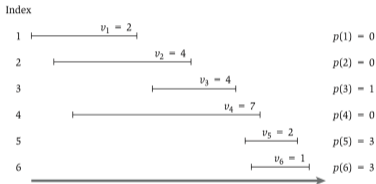


Figure 6.2 An instance of weighted interval scheduling with the functions $p(j)$ defined for each interval j .

E' 3 nella soluzione ottimale? SI

$$OPT(j) = \max \{ v_j + OPT(p(j)), OPT(j-1) \}$$

$$OPT(3) = \max \{ v_3 + OPT(p(3)), OPT(2) \} = \max \{ 4 + OPT(1), OPT(2) \} = \max \{ 4 + 2, 4 \}$$

$$M = \begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 1 & 2 & 3 & 4 & 5 & 6 \\ \hline 0 & 2 & 4 & 6 & 7 & 8 & 8 \\ \hline \end{array}$$

Schedulazione degli intervalli pesati: Trovare una soluzione

Domanda. Algoritmo di programmazione dinamica computa il valore ottimale. E se voglio una soluzione?

Risposta. Necessarie operazioni di post-processing.

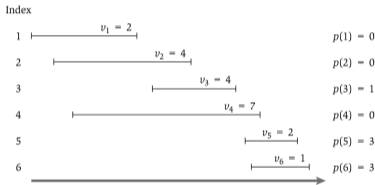


Figure 6.2 An instance of weighted interval scheduling with the functions $p(j)$ defined for each interval j .

E' 1 nella soluzione ottimale? SI

$$OPT(j) = \max \{ v_j + OPT(p(j)), OPT(j-1) \}$$

$$OPT(1) = \max \{ v_1 + OPT(p(1)), OPT(0) \} = \max \{ 2 + OPT(0), OPT(0) \} = \max \{ 2 + 0, 0 \}$$

$$M = \begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 1 & 2 & 3 & 4 & 5 & 6 \\ \hline 0 & 2 & 4 & 6 & 7 & 8 & 8 \\ \hline \end{array}$$

Schedulazione degli intervalli pesati: Trovare una soluzione

Domanda. Algoritmo di programmazione dinamica computa il valore ottimale. E se voglio una soluzione?

Risposta. Necessarie operazioni di post-processing.

```
Run M-Compute-Opt(n)
Run Find-Solution(n)

Find-Solution(j) {
  if (j = 0)
    output nothing
  else if ( $v_j + M[p(j)] > M[j-1]$ )
    print j
    Find-Solution(p(j))
  else
    Find-Solution(j-1)
}
```

Numero di chiamate ricorsive $\leq n \Rightarrow O(n)$.

Schedulazione degli intervalli pesati: Approccio bottom-Up

Programmazione dinamica bottom-up.

Input: $n, s_1, \dots, s_n, f_1, \dots, f_n, v_1, \dots, v_n$

Ordinare job per tempo di fine $f_1 \leq f_2 \leq \dots \leq f_n$.

Computa $p(1), p(2), \dots, p(n)$

```
Iterative-Compute-Opt {  
    M[0] = 0  
    for j = 1 to n  
        M[j] = max(v_j + M[p(j)], M[j-1])  
    }
```

Complessità $O(n)$.

Schedulazione degli intervalli pesati: Approccio bottom-Up

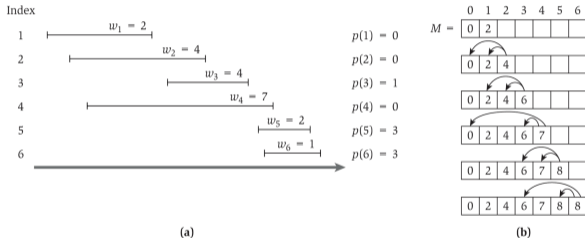


Figure 6.5 Part (b) shows the iterations of *Iterative-Compute-Opt* on the sample instance of Weighted Interval Scheduling depicted in part (a).

Esercizio: calcolo $p(j)$

Computazione di $p(\cdot)$: $O(n)$ dopo un ordinamento per tempo di inizio.

Quindi, dati i job ordinati per tempo di inizio, come computare in tempo lineare i valori $p(j)$?

Primo algoritmo:

```
Input:  $n, s_1, \dots, s_n, f_1, \dots, f_n$ 
```

```
job ordinati per tempo di fine  $f_1 \leq f_2 \leq \dots \leq f_n$ .
```

```
Computa-valori-p {
```

```
   $f_0=0$ 
```

```
  for  $i=1$  to  $n$ 
```

```
    Calcola il max  $j$  tale che  $f_j \leq s_i$ 
```

```
     $p(i)=j$ 
```

```
}
```


Esercizio: calcolo $p(j)$

Computazione di $p(\cdot)$: $O(n)$ dopo un ordinamento per tempo di inizio.

Quindi, dati i job ordinati per tempo di inizio, come computare in tempo lineare i valori $p(j)$?

Primo algoritmo:

```
Input:  $n, s_1, \dots, s_n, f_1, \dots, f_n$ 
```

```
job ordinati per tempo di fine  $f_1 \leq f_2 \leq \dots \leq f_n$ .
```

```
Computa-valori-p {
```

```
   $f_0=0$ 
```

```
  for  $i=1$  to  $n$ 
```

```
    Calcola il max  $j$  tale che  $f_j \leq s_i$ 
```

```
     $p(i)=j$ 
```

```
}
```

Comlessità?

Esercizio: calcolo $p(j)$

Computazione di $p(\cdot)$: $O(n)$ dopo un ordinamento per tempo di inizio.

Quindi, dati i job ordinati per tempo di inizio, come computare in tempo lineare i valori $p(j)$?

Input: $n, s_1, \dots, s_n, f_1, \dots, f_n, \pi$

job ordinati per tempo di fine $f_1 \leq f_2 \leq \dots \leq f_n$.

job ordinati per tempo di inizio $s_{\pi(1)} \leq s_{\pi(2)} \leq \dots \leq s_{\pi(n)}$.

```
Computa-valori-p {  
  j=0  
  f_0=0  
  for i=1 to n  
    Calcola il max j tale che  $f_j \leq s_{\pi(i)}$   
     $p(\pi(i))=j$   
}
```

Esercizio: calcolo $p(j)$

Computazione di $p(\cdot)$: $O(n)$ dopo un ordinamento per tempo di inizio.

Quindi, dati i job ordinati per tempo di inizio, come computare in tempo lineare i valori $p(j)$?

Input: $n, s_1, \dots, s_n, f_1, \dots, f_n, \pi$

job ordinati per tempo di fine $f_1 \leq f_2 \leq \dots \leq f_n$.

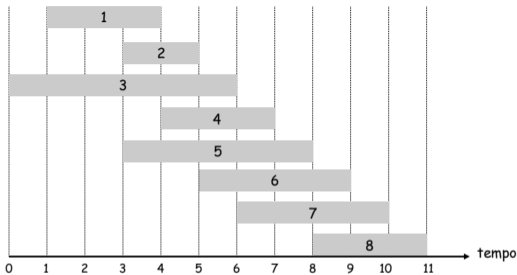
job ordinati per tempo di inizio $s_{\pi(1)} \leq s_{\pi(2)} \leq \dots \leq s_{\pi(n)}$.

```
Computa-valori-p {  
  j=0  
  f_0=0  
  for i=1 to n  
    Calcola il max j tale che  $f_j \leq s_{\pi(i)}$   
     $p(\pi(i))=j$   
}
```

$O(n)$

Non può essere più piccolo del valore j precedente!

Esercizio: calcolo $p(j)$, esempio



$s_{\pi(i)}$	0	1	3	3	4	5	6	8
$\pi(i)$	3	1	2	5	4	6	7	8

Iniziamo con $\pi(1)$

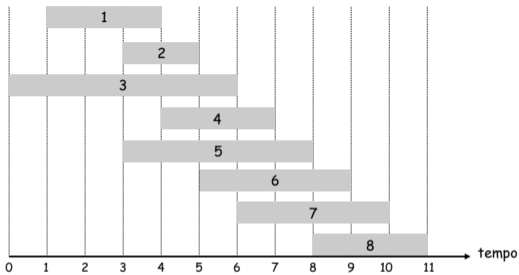
$s_{\pi(1)}=0$ ed $f_1=4$

Quindi $p(\pi(1))=0$, cioè $p(3)=0$

f_i	4	5	6	7	8	9	10	11
i	1	2	3	4	5	6	7	8

$P(\pi(i))$	0	0	0	0	1	2	3	5
$\pi(i)$	3	1	2	5	4	6	7	8

Esercizio: calcolo $p(j)$, esempio



$s_{\pi(i)}$	0	1	3	3	4	5	6	8
$\pi(i)$	3	1	2	5	4	6	7	8

Continuiamo con $\pi(2)$

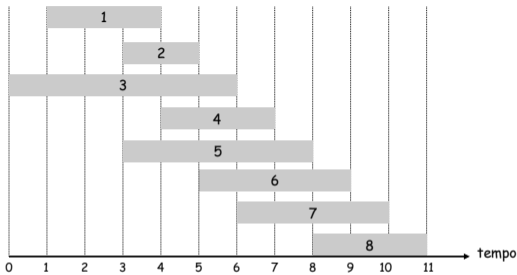
$s_{\pi(2)}=1$ ed $f_1=4$

Quindi $p(\pi(2))=0$, cioè $p(1)=0$

f_i	4	5	6	7	8	9	10	11
i	1	2	3	4	5	6	7	8

$P(\pi(i))$	0	0	0	0	1	2	3	5
$\pi(i)$	3	1	2	5	4	6	7	8

Esercizio: calcolo $p(j)$, esempio



$s_{\pi(i)}$	0	1	3	3	4	5	6	8
$\pi(i)$	3	1	2	5	4	6	7	8

Continuiamo con $\pi(3)$

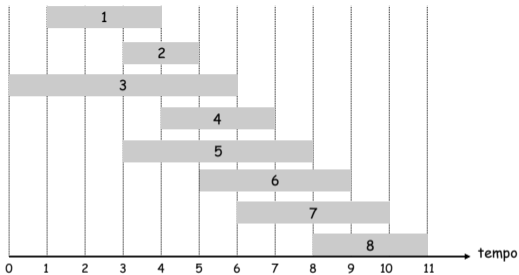
$s_{\pi(3)}=3$ ed $f_1=4$

Quindi $p(\pi(3))=0$, cioè $p(2)=0$

f_i	4	5	6	7	8	9	10	11
i	1	2	3	4	5	6	7	8

$P(\pi(i))$	0	0	0	0	1	2	3	5
$\pi(i)$	3	1	2	5	4	6	7	8

Esercizio: calcolo $p(j)$, esempio



$s_{\pi(i)}$	0	1	3	3	4	5	6	8
$\pi(i)$	3	1	2	5	4	6	7	8

Continuiamo con $\pi(4)$

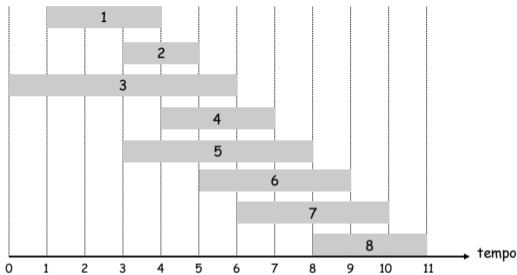
$s_{\pi(4)}=3$ ed $f_1=4$

Quindi $p(\pi(4))=0$, cioè $p(5)=0$

f_i	4	5	6	7	8	9	10	11
i	1	2	3	4	5	6	7	8

$P(\pi(i))$	0	0	0	0	1	2	3	5
$\pi(i)$	3	1	2	5	4	6	7	8

Esercizio: calcolo $p(j)$, esempio



$s_{\pi(i)}$	0	1	3	3	4	5	6	8
$\pi(i)$	3	1	2	5	4	6	7	8

Continuiamo con $\pi(5)$

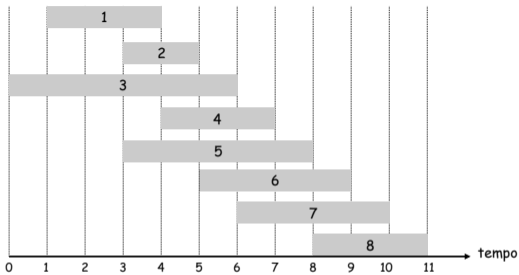
$s_{\pi(5)}=4$ ed $f_1=4, f_2=5$

Quindi $p(\pi(5))=1$, cioè $p(4)=1$

f_i	4	5	6	7	8	9	10	11
i	1	2	3	4	5	6	7	8

$P(\pi(i))$	0	0	0	0	1	2	3	5
$\pi(i)$	3	1	2	5	4	6	7	8

Esercizio: calcolo $p(j)$, esempio



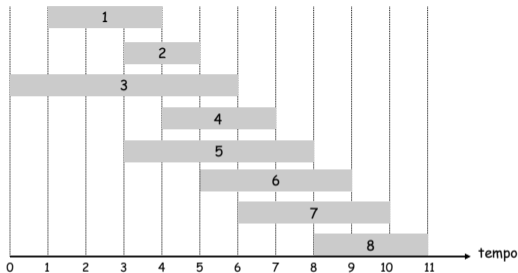
$s_{\pi(i)}$	0	1	3	3	4	5	6	8
$\pi(i)$	3	1	2	5	4	6	7	8

Continuiamo con $\pi(6)$
 $s_{\pi(6)}=5$ ed $f_1=4, f_2=5, f_3=6$
 Quindi $p(\pi(6))=1$, cioè $p(6)=2$

f_i	4	5	6	7	8	9	10	11
i	1	2	3	4	5	6	7	8

$P(\pi(i))$	0	0	0	0	1	2	3	5
$\pi(i)$	3	1	2	5	4	6	7	8

Esercizio: calcolo $p(j)$, esempio



$s_{\pi(i)}$	0	1	3	3	4	5	6	8
$\pi(i)$	3	1	2	5	4	6	7	8

Continuiamo con $\pi(7)$

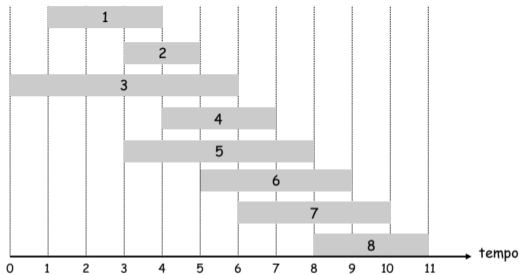
$s_{\pi(7)}=6$ ed $f_2=5, f_3=6, f_4=7$

Quindi $p(\pi(7))=3$, cioè $p(7)=3$

f_i	4	5	6	7	8	9	10	11
i	1	2	3	4	5	6	7	8

$P(\pi(i))$	0	0	0	0	1	2	3	5
$\pi(i)$	3	1	2	5	4	6	7	8

Esercizio: calcolo $p(j)$, esempio



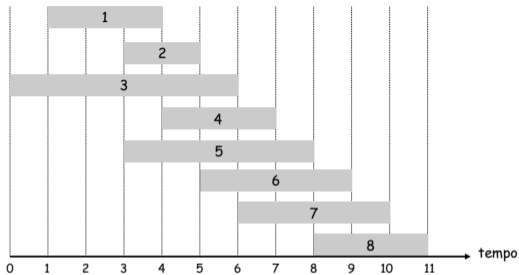
$s_{\pi(i)}$	0	1	3	3	4	5	6	8
$\pi(i)$	3	1	2	5	4	6	7	8

Terminiamo con $\pi(8)$
 $s_{\pi(8)}=8$ ed $f_3=6, f_4=7, f_5=8, f_6=9$
 Quindi $p(\pi(8))=5$, cioè $p(8)=5$

f_i	4	5	6	7	8	9	10	11
i	1	2	3	4	5	6	7	8

$P(\pi(i))$	0	0	0	0	1	2	3	5
$\pi(i)$	3	1	2	5	4	6	7	8

Esercizio: calcolo $p(j)$, esempio



$s_{\pi(i)}$	0	1	3	3	4	5	6	8
$\pi(i)$	3	1	2	5	4	6	7	8

$P(i)$	0	0	0	1	0	2	3	5
i	1	2	3	4	5	6	7	8

f_i	4	5	6	7	8	9	10	11
i	1	2	3	4	5	6	7	8

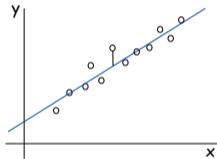
$P(\pi(i))$	0	0	0	0	1	2	3	5
$\pi(i)$	3	1	2	5	4	6	7	8

6.3 Segmented Least Squares

Segmented Least Squares

Least squares.

- Problema fondamentale in statistica ed analisi numerica.
- Dati n punti nel piano: $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$.
- Trovare linea $y = ax + b$ che minimizza la somma dell'errore quadratico:

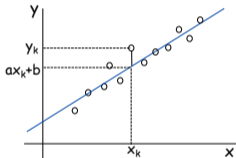


Segmented Least Squares

Least squares.

- Problema fondamentale in statistica ed analisi numerica.
- Dati n punti nel piano: $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$.
- Trovare linea $y = ax + b$ che minimizza la somma dell'errore quadratico:

$$\text{Errore}("y = ax + b", \text{punti}) = \sum_{k=1}^n (y_k - ax_k - b)^2$$

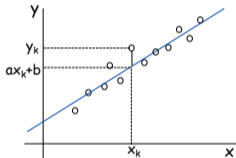


Segmented Least Squares

Least squares.

- Problema fondamentale in statistica ed analisi numerica.
- Dati n punti nel piano: $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$.
- Trovare linea $y = ax + b$ che minimizza la somma dell'errore quadratico:

$$\text{Errore}("y = ax + b", \text{punti}) = \sum_{k=1}^n (y_k - ax_k - b)^2$$



Soluzione. Sappiamo che l'errore minimo si ha quando

$$a = \frac{n \sum_k x_k y_k - (\sum_k x_k) (\sum_k y_k)}{n \sum_k x_k^2 - (\sum_k x_k)^2}, \quad b = \frac{\sum_k y_k - a \sum_k x_k}{n}$$

Insieme di punti su due linee

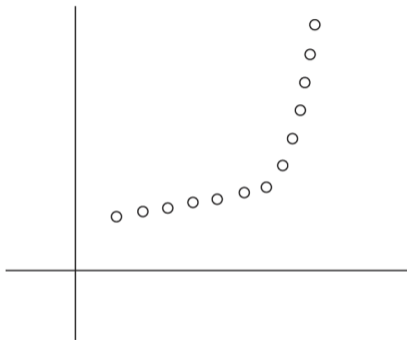


Figure 6.7 A set of points that lie approximately on two lines.

Insieme di punti su tre linee

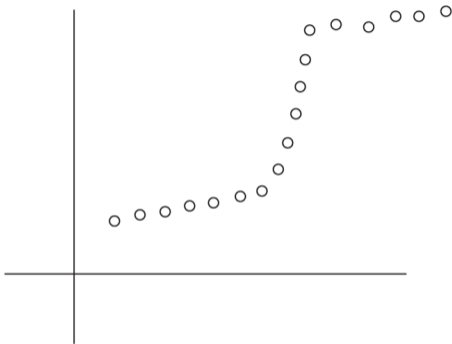


Figure 6.8 A set of points that lie approximately on three lines.

Segmented Least Squares

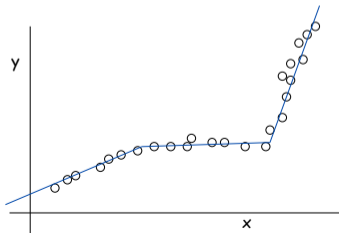
Segmented least squares.

- Punti si trovano approssimativamente su una sequenza di segmenti.
- Dati n punti nel piano $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ con $x_1 < x_2 < \dots < x_n$, trovare una sequenza di linee che minimizza $f(x)$.

Domanda. Qual'è una scelta ragionevole per $f(x)$ per bilanciare accuratezza e parsominia?

↑
numero di linee

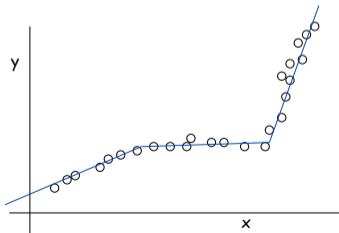
↑
bontà dell'approssimazione



Segmented Least Squares

Segmented least squares.

- Punti si trovano approssimativamente su una sequenza di segmenti.
- Dati n punti nel piano $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ con $x_1 < x_2 < \dots < x_n$, trovare una sequenza di linee che minimizza :
 - la somma delle somme degli errori quadratici E in ogni segmento
 - il numero delle linee L
- Funzione tradeoff: $E + c L$, per qualche costante $c > 0$.



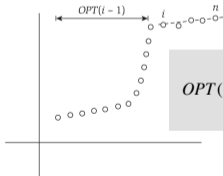
Programmazione dinamica: Multiway Choice

Notazione.

- $OPT(j)$ = minimo costo per punti p_1, p_2, \dots, p_j . [con $OPT(0)=0$]
- $e(i, j)$ = minima somma di errori quadratici per punti p_i, p_{i+1}, \dots, p_j .

Per computare $OPT(j)$:

- L'ultimo segmento usa punti p_i, p_{i+1}, \dots, p_j per qualche i .
- Costo = $e(i, j) + c + OPT(i-1)$.



$$OPT(j) = \begin{cases} 0 & \text{se } j = 0 \\ \min_{1 \leq i \leq j} \{ e(i, j) + c + OPT(i-1) \} & \text{altrimenti} \end{cases}$$

Figure 6.9 A possible solution: a single line segment fits points p_i, p_{i+1}, \dots, p_n , and then an optimal solution is found for the remaining points p_1, p_2, \dots, p_{i-1} .

Segmented Least Squares: Algoritmo

```
INPUT:  $n, p_1, \dots, p_N, c$ 

Segmented-Least-Squares() {
  M[0] = 0
  for j = 1 to n
    for i = 1 to j
      compute il minimo errore quadratico  $e_{ij}$ 
      per il segmento  $p_1, \dots, p_j$ 

  for j = 1 to n
    M[j] =  $\min_{1 \leq i \leq j} (e_{ij} + c + M[i-1])$ 

  return M[n]
}
```

Complessità. $O(n^3)$.

- Collo di bottiglia = computazione di $e(i, j)$ per $O(n^2)$ coppie, $O(n)$ per coppia usando la formula precedente

Segmented Least Squares: Algoritmo

```
INPUT:  $n, p_1, \dots, p_N, c$ 

Segmented-Least-Squares() {
  M[0] = 0
  for j = 1 to n
    for i = 1 to j
      compute the least square error  $e_{ij}$  for
      the segment  $p_i, \dots, p_j$ 

  for j = 1 to n
    M[j] =  $\min_{1 \leq i \leq j} (e_{ij} + c + M[i-1])$ 

  return M[n]
}
```

Complessità: $O(n^3)$

può essere migliorato: $O(n^2)$ pre-computando varie statistiche
Esercizio da fare a casa

Idea: calcolare valori (i,j) per $j-i=1, j-i=2, \dots$
valori (i,j) da valori $(i,j-1)$ in $O(1)$

Collo di bottiglia = computazione di $e(i, j)$ per $O(n^2)$ coppie, $O(n)$ per coppia usando la formula precedente

Esercizio: calcolare tutti i valori $e(i,j)$ in $O(n^2)$

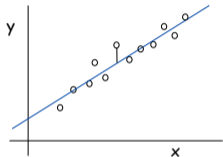
Algoritmo di forza bruta: $O(n^3)$

Suggerimento per algoritmo $O(n^2)$:

- ❑ calcolare valori $e(i,j)$ per $j-i=1, j-i=2, \dots$
- ❑ calcolare $e(i,j)$ da $e(i,j-1)$ in $O(1)$

$$e(i,j) = \sum_{k=i}^j (y_k - ax_k - b)^2$$

$$a = \frac{n \sum_k x_k y_k - (\sum_k x_k)(\sum_k y_k)}{n \sum_k x_k^2 - (\sum_k x_k)^2}, \quad b = \frac{\sum_k y_k - a \sum_k x_k}{n}$$



Esercizio: calcolare tutti i valori $e(i,j)$ in $O(n^2)$

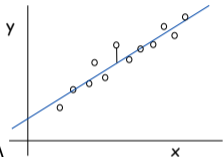
Algoritmo di forza bruta: $O(n^3)$

Suggerimento per algoritmo $O(n^2)$:

- ❑ calcolare valori (i,j) per $j-i=1, j-i=2, \dots$
- ❑ calcolare valori (i,j) da quelli $(i,j-1)$ in $O(1)$

$$e(i,j) = \sum_{k=i}^j (y_k - ax_k - b)^2$$

$$a = \frac{n \sum_k x_k y_k - (\sum_k x_k)(\sum_k y_k)}{n \sum_k x_k^2 - (\sum_k x_k)^2}, \quad b = \frac{\sum_k y_k - a \sum_k x_k}{n}$$



Soluzione: calcolare tutte le somme

$$\sum_{k=i}^j x_k, \quad \sum_{k=i}^j y_k, \quad \sum_{k=i}^j x_k y_k, \quad \sum_{k=i}^j x_k^2, \quad \sum_{k=i}^j y_k^2 \quad \text{in } O(n^2)$$

Scrivere l'errore come

$$\sum_{k=i}^j (y_k - ax_k - b)^2 = \sum_{k=i}^j y_k^2 + a^2 \sum_{k=i}^j x_k^2 + b^2 n - 2a \sum_{k=i}^j x_k y_k - 2b \sum_{k=i}^j y_k + 2ab \sum_{k=i}^j x_k$$

Segmented Least Squares: Trovare una soluzione

Find-Segments(j)

 If $j=0$ then

 Output nothing

 Else

 Find an i that minimizes $e_{i,j} + C + M[i-1]$

 Output the segment $\{p_i, \dots, p_j\}$ and the result of

 Find-Segments($i-1$)

 Endif

6.4 Knapsack Problem

Problema dello zaino

Problema dello zaino.

- Abbiamo n oggetti ed uno "zaino"
- Oggetto i pesa $w_i > 0$ chilogrammi ed ha valore $v_i > 0$.
- Zaino ha una capacità di W chilogrammi.
- Obiettivo: riempire zaino per massimizzare valore totale.

Esempio: { 3, 4 } ha valore 40.

$$W = 11$$

oggetti	valore	peso
1	1	1
2	6	2
3	18	5
4	22	6
5	28	7

Algoritmo greedy: aggiungere oggetto con peso massimo w_i .

Esempio: { 5, 2, 1 } ha valore = 35 \Rightarrow algoritmo greedy non ottimale.

Problema dello zaino

Problema dello zaino.

- Abbiamo n oggetti ed uno "zaino."
- Oggetto i pesa $w_i > 0$ chilogrammi ed ha valore $v_i > 0$.
- Zaino ha una capacità di W chilogrammi.
- Obiettivo: riempire zaino per massimizzare valore totale.

Esempio: { 3, 4 } ha valore 40.

$$W = 11$$

oggetti	valore	peso
1	1	1
2	6	2
3	18	5
4	22	6
5	28	7

Algoritmo greedy: aggiungere oggetto con valore massimo v_i .

Esempio: { 5, 2, 1 } ha valore = 35 \Rightarrow algoritmo greedy non ottimale.

Problema dello zaino

Problema dello zaino.

- Abbiamo n oggetti ed uno "zaino."
- Oggetto i pesa $w_i > 0$ chilogrammi ed ha valore $v_i > 0$.
- Zaino ha una capacità di W chilogrammi.
- Obiettivo: riempire zaino per massimizzare valore totale.

Esempio: { 3, 4 } ha valore 40.

$$W = 11$$

oggetti	valore	peso	val/peso
1	1	1	1
2	6	2	3
3	18	5	3,6
4	22	6	3,66
5	28	7	4

Algoritmo greedy: aggiungere oggetto con rapporto massimo v_i / w_i

Esempio: { 5, 2, 1 } ha valore = 35 \Rightarrow algoritmo greedy non ottimale.

Programmazione dinamica: Primo tentativo

Definizione. $OPT(i)$ = max valore totale con sottoinsieme oggetti 1, ..., i.

- **Caso 1:** OPT non contiene oggetto i .
 - OPT contiene max valore totale di $\{ 1, 2, \dots, i-1 \}$
- **Caso 2:** OPT contiene oggetto i .
 - che l'oggetto i ci sia nella soluzione non implica immediatamente che altri oggetti non ci siano
 - senza sapere quali altri oggetti ci siano nella soluzione non sappiamo se c'è spazio sufficiente per altri oggetti

Conclusione. Abbiamo bisogno di più sottoproblemi!

Programmazione dinamica: aggiunta nuova variabile

Definizione. $OPT(i, w) = \max$ valore totale con sottoinsieme oggetti $1, \dots, i$ con limite di peso w .

- Case 1: OPT non contiene oggetto i .
 - OPT contiene max valore totale di $\{ 1, 2, \dots, i-1 \}$ con limite di peso w
- Case 2: OPT contiene oggetto i .
 - Nuovo limite di peso = $w - w_i$
 - OPT contiene max valore totale di $\{ 1, 2, \dots, i-1 \}$ con nuovo limite peso

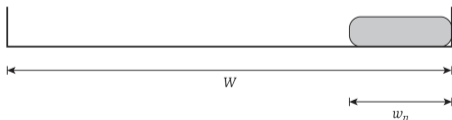


Figure 6.10 After item n is included in the solution, a weight of w_n is used up and there is $W - w_n$ available weight left.

Programmazione dinamica: aggiunta nuova variabile

Definizione. $OPT(i, w) = \max$ valore totale con sottoinsieme oggetti $1, \dots, i$ con limite di peso w .

- Case 1: OPT non contiene oggetto i .
 - OPT contiene max valore totale di $\{1, 2, \dots, i-1\}$ con limite di peso w
- Case 2: OPT contiene oggetto i .
 - Nuovo limite di peso = $w - w_i$
 - OPT contiene max valore totale di $\{1, 2, \dots, i-1\}$ con nuovo limite di peso

$$OPT(i, w) = \begin{cases} 0 & \text{se } i = 0 \\ OPT(i-1, w) & \text{se } w_i > w \\ \max\{OPT(i-1, w), v_i + OPT(i-1, w - w_i)\} & \text{altrimenti} \end{cases}$$

Problema dello zaino: approccio bottom-up

Utilizzare un array n -per- W .

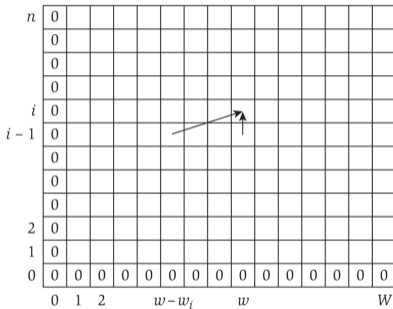


Figure 6.11 The two-dimensional table of OPT values. The leftmost column and bottom row is always 0. The entry for $OPT(i, w)$ is computed from the two other entries $OPT(i - 1, w)$ and $OPT(i - 1, w - w_i)$, as indicated by the arrows.

Problema dello zaino: approccio bottom-up

Utilizzare un array n-per-W.

$$OPT(i, w) = \begin{cases} 0 & \text{se } i=0 \\ OPT(i-1, w) & \text{se } w_i > w \\ \max\{OPT(i-1, w), v_i + OPT(i-1, w-w_i)\} & \text{altrimenti} \end{cases}$$

```
Input: n, w1, ..., wN, v1, ..., vN
```

```
for w = 0 to W
```

```
  M[0, w] = 0
```

```
for i = 1 to n
```

```
  for w = 1 to W
```

```
    if (wi > w)
```

```
      M[i, w] = M[i-1, w]
```

```
    else
```

```
      M[i, w] = max {M[i-1, w], vi + M[i-1, w-wi ]}
```

```
return M[n, W]
```

Problema dello zaino: esempio algoritmo

Knapsack size $W = 6$, items $v_1 = 2, v_2 = 2, v_3 = 1$
 $w_1 = 2, w_2 = 2, w_3 = 3$

3							
2							
1							
0	0	0	0	0	0	0	0
	0	1	2	3	4	5	6

Initial values

3							
2							
①	0	0	2	2	2	2	2
0	0	0	0	0	0	0	0
	0	1	2	3	4	5	6

Filling in values for $i = 1$

3							
②	0	0	2	2	4	4	4
1	0	0	2	2	2	2	2
0	0	0	0	0	0	0	0
	0	1	2	3	4	5	6

Filling in values for $i = 2$

3							
③	0	0	2	3	4	5	5
2	0	0	2	2	4	4	4
1	0	0	2	2	2	2	2
0	0	0	0	0	0	0	0
	0	1	2	3	4	5	6

Filling in values for $i = 3$

Figure 6.12 The iterations of the algorithm on a sample instance of the Subset Sum Problem.

Problema dello zaino: esempio algoritmo

		$\xrightarrow{\hspace{10em}} W + 1 \xrightarrow{\hspace{10em}}$											
		0	1	2	3	4	5	6	7	8	9	10	11
$n + 1$	ϕ	0	0	0	0	0	0	0	0	0	0	0	0
	{1}	0	1	1	1	1	1	1	1	1	1	1	1
	{1, 2}	0	1	6	7	7	7	7	7	7	7	7	7
	{1, 2, 3}	0	1	6	7	7	18	19	24	25	25	25	25
	{1, 2, 3, 4}	0	1	6	7	7	18	22	24	28	29	29	40
	{1, 2, 3, 4, 5}	0	1	6	7	7	18	22	28	29	34	34	40

OPT: {4, 3}
 valore = 22 + 18 = 40

W = 11

oggetto	valore	peso
1	1	1
2	6	2
3	18	5
4	22	6
5	28	7

Problema dello zaino: Complessità

Complessità. $\Theta(nW)$.

- Non polinomiale nella lunghezza dell' input!
- "Pseudo-polinomiale."
- Versione decisionale del problema dello zaino è NP-completo. [Cap 8]

Algoritmi di approssimazione per lo zaino. Esiste un algoritmo di approssimazione polinomiale che produce una soluzione il cui valore è al massimo 0.01% rispetto all' ottimale. [Sezione 11.8]

Comunque, non lo tratteremo in questo corso!

Problema dello zaino: Esercizio 1

Si descriva ed analizzi un algoritmo per la seguente variazione del problema dello zaino: Dati n oggetti di peso w_1, w_2, \dots, w_n e valore v_1, v_2, \dots, v_n ed uno zaino di capacità W (tutti gli input sono >0), trovare il massimo valore di un sottoinsieme degli oggetti il cui peso totale è al massimo W , con la condizione che ogni oggetto può essere preso anche più di una volta.

(La variazione rispetto al problema del testo, consiste nel superamento del vincolo che ogni oggetto poteva essere preso al massimo una sola volta.)

Problema dello zaino: Esercizio 2

Si descriva ed analizzi un algoritmo per la seguente variazione del problema dello zaino: Dati n oggetti di peso w_1, w_2, \dots, w_n e valore v_1, v_2, \dots, v_n ed uno zaino di capacità W (tutti gli input sono >0), trovare il massimo valore di un sottoinsieme degli oggetti il cui peso totale è al massimo W , con la condizione che ogni oggetto può essere preso al massimo 2 volte.

(La variazione rispetto al problema del testo, consiste nel superamento del vincolo che ogni oggetto poteva essere preso al massimo una sola volta.)

Problema dello zaino: Esercizio 3

Si descriva ed analizzi un algoritmo per la seguente variazione del problema dello zaino: Dati n oggetti di peso w_1, w_2, \dots, w_n e valore v_1, v_2, \dots, v_n ed uno zaino di capacità W (tutti gli input sono >0), trovare il massimo valore di un sottoinsieme degli oggetti il cui peso totale è al massimo W , con la condizione che non possono essere presi due oggetti con indici consecutivi (ovvero gli oggetti i -esimo ed $(i+1)$ -esimo, per $i=1, 2, \dots, n-1$).


6.5 RNA Secondary Structure

Human Genome Project

Iniziato Ottobre 1990

Durata prevista 15 anni

Scopo: determinare la sequenza completa delle basi (3 miliardi) del DNA, identificare tutti i geni umani e renderli accessibili per ulteriori studi



Human Genome Project

Impacting many disciplines

*Courtesy
U.S. Department of Energy
Human Genome Program*

*Global Carbon Cycles
Industrial Resources • Bioremediation
Evolutionary Biology • Biofuels • Agriculture • Forensics
Molecular and Nuclear Medicine • Health Risks*

HGA 06 11108

Genoma Umano

Fatto di DNA che ha 4 differenti mattoni chimici, chiamati basi: A, T, C, G.

Se la sequenza fosse scritta in un elenco telefonico, occorrerebbero 200 volumi di 100 pagine ognuno

Se provassimo a leggere 10 basi al secondo, cioè 600 basi/minuto, 36.000 basi/ora, 864.000 basi/giorno, 315.360.000 basi/anno, occorrerebbero circa 9,5 anni per leggere la sequenza
1 Mb = 1.000.000 basi (megabase)

Per intera sequenza occorrono 3 gigabyte, senza contare altri dati associati

Genoma Umano

Human Genome Research Institute (finanziamenti pubblici USA) e
Celera Genomics Corporation (settore commerciale)

Sequenze geniche sono brevettabili, stabilito dai tribunali (brevetto
dura 20 anni)

Corsa ai brevetti per Celera ed alla pubblicazione degli accademici per
evitare i brevetti

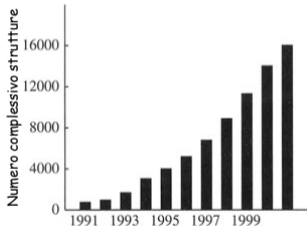
Corsa finita alla pari

26 giugno 2000, Tony Blair e Bill Clinton
annunciano il sequenziamento completo
del genoma umano

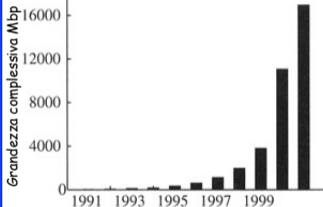


Organism	Genome Size (Bases)	Estimated Genes
Human (<i>Homo sapiens</i>)	3 billion	30,000
Laboratory mouse (<i>M. musculus</i>)	2.6 billion	30,000
Mustard weed (<i>A. thaliana</i>)	100 million	25,000
Roundworm (<i>C. elegans</i>)	97 million	19,000
Fruit fly (<i>D. melanogaster</i>)	137 million	13,000
Yeast (<i>S. cerevisiae</i>)	12.1 million	6,000
Bacterium (<i>E. coli</i>)	4.6 million	3,200
Human immunodeficiency virus (HIV)	9700	9

Crescita Banche Dati



Protein Data Bank, archivio di strutture tridimensionali di macromolecole biologiche



GeneBank, banca dati di sequenze geniche del *National Center for Biotechnology Information* degli USA

Ogni organismo ha un genoma che contiene tutte le informazioni biologiche necessarie per costruire e mantenere un esempio vivente di quell'organismo

Cromosoma: composto da una molecola di Dna e dalle proteine associate, contiene informazioni ereditare

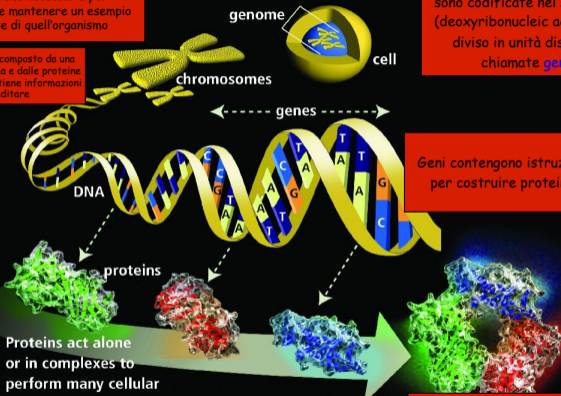
Le informazioni biologiche sono codificate nel suo DNA (deoxyribonucleic acid) ed è diviso in unità discrete chiamate geni

Geni contengono istruzioni per costruire proteine

Proteins act alone or in complexes to perform many cellular functions

Le proteine determinano come appare l'organismo, come il corpo metabolizza cibo o combatte infezioni, e qualche volta come si comporta, ...

U.S. DEPARTMENT OF ENERGY



Genoma

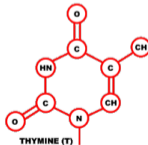
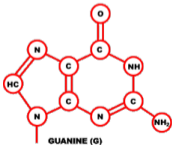
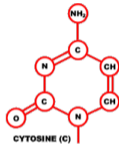
Ogni organismo ha un genoma che contiene tutte le informazioni biologiche necessarie per costruire mantenere un esempio vivente di quell'organismo

Le informazioni biologiche sono codificate nel suo DNA (deoxyribonucleic acid) ed è diviso in unità discrete chiamate geni

Geni contengono istruzioni per costruire proteine richieste dall'organismo

Le proteine determinano, come appare l'organismo, come il corpo metabolizza cibo o combatte infezioni, e qualche volta come si comporta, ...

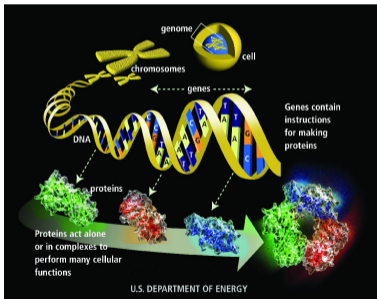
Le 4 basi del DNA



DNA

Molecola formata da 2 catene a forma di doppia elica

Catene connesse da basi complementari A-T e C-G



equipment, and to Dr. G. E. R. Denton and the origin and address of R.R.E. *Discovery II* for their part in making the observations.

*Tovey, P. R., General, R., and James, W., *Phil. Mag.*, **46**, 133 (1953).

*Langan Taylor, M. S., *Ann. Ent. Soc. Amer.*, **46**, 663 (1953).

*Van Arman, W., *Woods Hole Paper in Phys. Chem. Botany*, **22**, 11 (1952).

*Simons, V. W., *Ann. Ent. Soc. Amer.*, **46**, 111 (1953).

MOLECULAR STRUCTURE OF NUCLEIC ACIDS

A Structure for Deoxyribose Nucleic Acid

WE wish to suggest a structure for the salt of deoxyribose nucleic acid (D.N.A.). This structure has novel features which are of considerable biological interest.

A structure for nucleic acid has already been proposed by Pauling and Corey¹. They kindly made their manuscript available to us in advance of publication. Their model consists of three inter-twined chains, with the phosphates near the fibre axis, and the bases on the outside. In our opinion, this structure is unsatisfactory for two reasons:

(1) We believe that the material which gives the X-ray diagrams is the salt, not the free acid. Without the acidic hydrogen atoms it is not clear what forces would hold the structure together, especially as the negatively charged phosphates near the axis will repel each other. (2) Some of the van der Waals distances appear to be too small.

Another three-chain structure has also been suggested by Fraser (in the press). In his model the phosphates are on the outside and the bases on the inside, linked together by hydrogen bonds. This structure as described in neither is defined, and for this reason we shall not comment on it.

We wish to put forward a radically different structure for the salt of deoxyribose nucleic acid. This structure has two helical chains each coiled round the same axis (see diagram). We have made the usual chemical analysis, and find that the chains consist of phosphate di-ester groups joining 5'-deoxy-ribose nucleoside with 3',5' triphosphates. The two chains (but not their bases) are related by a dyad perpendicular to the fibre axis. Both chains follow right-handed helices, but owing to the dyad the sequence of the atoms in the two chains run in opposite directions. Each chain loosely resembles Furberg's model No. 1; that is, the bases are on the inside of the helix and the phosphates on the outside. The configuration of the sugar and the atoms near it is close to Furberg's standard configuration, the sugar being roughly perpendicular to the attached base. There

This figure is a projection of the model. The two chains revolve in opposite directions about the same axis, and the bases are on the inside of the chains. The vertical line marks the fibre axis.

is a residue on each chain every 2.4 Å. in the z-direction. We have measured an angle of 36° between adjacent residues in the same chain, so that the structure repeats after 10 residues on each chain, that is, after 24 Å. The distance of a phosphate atom from the fibre axis is 10 Å. As the phosphates are on the outside, cations have easy access to them.

The structure is an open one, and its water content is rather high. At lower water contents we would expect the bases to tilt so that the structure would become more compact.

The novel features of the structure is the manner in which the two chains are held together by the purine and pyrimidine bases. The planes of the bases are perpendicular to the fibre axis. They are joined together in pairs, a single base from one chain being hydrogen-bonded to a single base from the other chain, so that the two lie side by side with identical z-co-ordinates. One of the pair must be a purine and the other a pyrimidine for bonding to occur. The hydrogen bonds are made as follows: purine position 1 to pyrimidine position 1; purine position 6 to pyrimidine position 4.

If it is assumed that the bases only occur in the structure in the most plausible tautomeric form (that is, with the lone rather than the end configuration) it is found that only specific pairs of bases can bond together. These pairs are adenine (purine) with thymine (pyrimidine), and guanine (purine) with cytosine (pyrimidine).

In other words, if an adenine forms one member of a pair on either chain, then on those assumptions the other member must be thymine; similarly for guanine and cytosine. The sequence of bases on a single chain does not appear to be restricted in any way. However, if only specific pairs of bases can be formed, it follows that if the sequence of bases on one chain is given, then the sequence on the other chain is automatically determined.

It has been found experimentally^{2,3} that the ratio of the amounts of adenine to thymine, and the ratio of guanine to cytosine, are always very close to unity for deoxyribose nucleic acid.

It is probably impossible to build this structure with a ribose sugar in place of the deoxyribose, as the extra oxygen atom would make too close a van der Waals contact.

The previously published X-ray data^{4,5} on deoxyribose nucleic acid are insufficient for a rigorous test of our structure. So far as we can tell, it is roughly compatible with the experimental data, but it must be regarded as unproved until it has been checked against more exact results. Some of these are given in the following communications. We were not aware of the details of the results presented there when we devised our structure, which rests mainly though not entirely on published experimental data and stereochemical arguments.

It has not escaped our notice that the specific pairing we have postulated immediately suggests a possible copying mechanism for the genetic material.

Full details of the structure, including the conditions assumed in building it, together with a set of co-ordinates for the atoms, will be published elsewhere.

We are much indebted to Dr. Jerry Donohue for constant advice and criticism, especially on inter-atomic distances. We have also been stimulated by a knowledge of the general nature of the unpublished experimental results and ideas of Dr. M. H. F. Wilkins, Dr. R. E. Franklin and their co-workers at

King's College, London. One of us (J. D. W.) has been aided by a fellowship from the National Foundation for Infectious Diseases.

J. D. WATSON
F. H. C. CRICK

Medical Research Council Unit for the Study of the Molecular Structure of Biological Systems, Cavendish Laboratory, Cambridge.

April 2

¹Pauling, L., and Corey, R. B., *Nature*, **157**, 305 (1945); *Proc. U.S. Nat. Acad. Sci.*, **30**, 41 (1945).

²Furberg, B., *Acta Chem. Scand.*, **4**, 436 (1950).

³Chapdel, R., for references see Zavadoff, S., Zavadoff, G., and Champel, R., *Biochim. Biophys. Acta*, **8**, 277 (1952).

⁴Wyatt, G. G., *J. Gen. Physiol.*, **34**, 305 (1952).

⁵Arthur, W. T., *Temp. Soc. Rep. Ser.*, **1**, 194 (1952) (*Nature*, **170**, 789, 1952).

⁶Wilkins, M. H. F., and Randall, J. T., *Biophys. J.*, **1**, 102 (1952).



James Watson e
Francis Crick

Ribonucleic acid (RNA)

Simile al DNA.

Singola catena con 4 nucleotidi: adenine (A), cytosine (C), guanine (G), uracil (U).

RNA. Stringa $B = b_1b_2\dots b_n$ su alfabeto $\{A, C, G, U\}$.

RNA Secondary Structure

RNA. Stringa $B = b_1b_2\dots b_n$ su alfabeto $\{A, C, G, U\}$.

Struttura secondaria. RNA è una singola catena e tende a formare coppie di basi con se stessa. Questa struttura è essenziale per capire il comportamento delle molecole.

coppie di base complementari: A-U, C-G

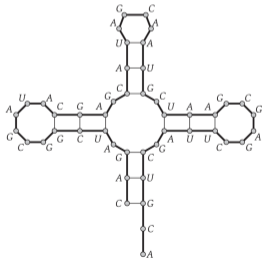


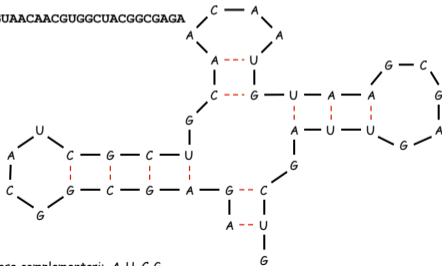
Figure 6.13 An RNA secondary structure. Thick lines connect adjacent elements of the sequence; thin lines indicate pairs of elements that are matched.

RNA Secondary Structure

RNA. Stringa $B = b_1b_2\dots b_n$ su alfabeto $\{A, C, G, U\}$.

Struttura secondaria. RNA è una singola catena e tende a formare coppie di basi con se stessa. Questa struttura è essenziale per capire il comportamento delle molecole.

Esempio: **GUCGAUUGAGCGAAUGUAACAACGUGGCUACGGCGAGA**



RNA Secondary Structure

Struttura secondaria. Un insieme di coppie $S = \{ (b_i, b_j) \}$ che soddisfa:

- [Watson-Crick.] S è un matching ed ogni coppia in S è un complemento Watson-Crick: A-U, U-A, C-G, oppure G-C.

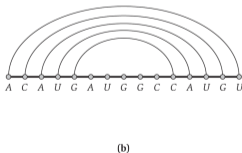
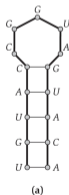
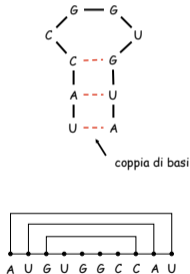


Figure 6.14 Two views of an RNA secondary structure. In the second view, (b), the string has been "stretched" lengthwise, and edges connecting matched pairs appear as noncrossing "bubbles" over the string.

RNA Secondary Structure

Struttura secondaria. Un insieme di coppie $S = \{ (b_i, b_j) \}$ che soddisfa:

- [Watson-Crick.] S è un matching ed ogni coppia in S è un complemento Watson-Crick: A-U, U-A, C-G, oppure G-C.
- [No sharp turns.] Gli elementi di ogni coppia sono separati da almeno 4 basi. Se $(b_i, b_j) \in S$, allora $i < j - 4$.

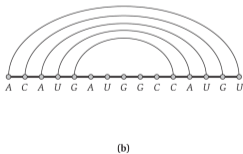
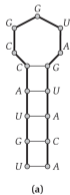
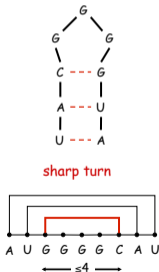


Figure 6.14 Two views of an RNA secondary structure. In the second view, (b), the string has been "stretched" lengthwise, and edges connecting matched pairs appear as noncrossing "bubbles" over the string.

RNA Secondary Structure

Struttura secondaria. Un insieme di coppie $S = \{ (b_i, b_j) \}$ che soddisfa:

- [Watson-Crick.] S è un matching ed ogni coppia in S è un complemento Watson-Crick: A-U, U-A, C-G, oppure G-C.
- [No sharp turns.] Gli elementi di ogni coppia sono separati da almeno 4 basi. Se $(b_i, b_j) \in S$, allora $i < j - 4$.
- [Non-crossing.] Se (b_i, b_j) e (b_k, b_l) sono due coppie in S , allora non è possibile che $i < k < j < l$.

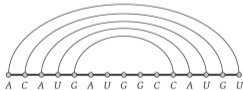
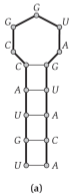
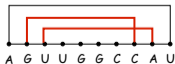
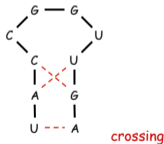


Figure 6.14 Two views of an RNA secondary structure. In the second view, (b), the string has been "stretched" lengthwise, and edges connecting matched pairs appear as noncrossing "bubbles" over the string.

RNA Secondary Structure

Struttura secondaria. Un insieme di coppie $S = \{ (b_i, b_j) \}$ che soddisfa:

- [Watson-Crick.] S è un matching ed ogni coppia in S è un complemento Watson-Crick: $A-U$, $U-A$, $C-G$, oppure $G-C$.
- [No sharp turns.] Gli elementi di ogni coppia sono separati da almeno 4 basi. Se $(b_i, b_j) \in S$, allora $i < j - 4$.
- [Non-crossing.] Se (b_i, b_j) e (b_k, b_l) sono due coppie in S , allora non è possibile che $i < k < j < l$.

Energia libera. L'ipotesi usuale è che una molecola RNA formerà la struttura secondaria con l'energia libera totale ottimale.

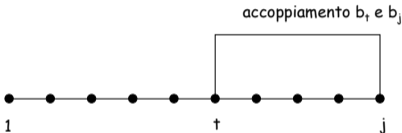
\

approssimativamente il numero di coppie di basi

Obiettivo. Data una molecola RNA $B = b_1b_2\dots b_n$, trovare una struttura secondaria S che massimizza il numero di coppie di basi.

RNA Secondary Structure: Sottoproblemi

Primo tentativo. $OPT(j)$ = massimo numero di coppie di basi in una struttura secondaria della sottostringa $b_1b_2\dots b_j$.



Difficoltà. Abbiamo due sottoproblemi:

- Trovare la struttura secondaria in: $b_1b_2\dots b_{t-1}$. ← $OPT(t-1)$
- Trovare la struttura secondaria in: $b_{t+1}b_{t+2}\dots b_{j-1}$. ← abbiamo bisogno di più sottoproblemi

Programmazione dinamica su intervalli

Notazione. $OPT(i, j)$ = massimo numero di coppie di basi in una struttura secondaria della sottostringa $b_i b_{i+1} \dots b_j$.

- **Case 1.** Se $i \geq j - 4$.
 - $OPT(i, j) = 0$ per la condizione “no-sharp turns”.
- **Case 2.** La base b_j non fa parte di una coppia.
 - $OPT(i, j) = OPT(i, j-1)$
- **Case 3.** La base b_j fa coppia con b_t per qualche $i \leq t < j - 4$.
 - Vincolo “non-crossing” determina i risultanti sottoproblemi
 - $OPT(i, j) = \max_{i \leq t < j-4} \{ 1 + OPT(i, t-1) + OPT(t+1, j-1) \}$

b_t e b_j sono complementi Watson-Crick, cioè A-U, U-A, C-G, oppure G-C.

Programmazione dinamica su intervalli

Notazione. $OPT(i, j)$ = massimo numero di coppie di basi in una struttura secondaria della sottostringa $b_i b_{i+1} \dots b_j$.

- **Case 1.** Se $i \geq j - 4$.
 - $OPT(i, j) = 0$ per la condizione “no-sharp turns”.
- **Case 2.** La base b_j non fa parte di una coppia.
 - $OPT(i, j) = OPT(i, j-1)$
- **Case 3.** La base b_j fa coppia con b_t per qualche $i \leq t < j - 4$.
 - Vincolo “non-crossing” determina i risultanti sottoproblemi
 - $OPT(i, j) = \max_{i \leq t < j-4} \{ 1 + OPT(i, t-1) + OPT(t+1, j-1) \}$

b_t e b_j sono complementi Watson-Crick, cioè A-U, U-A, C-G, oppure G-C.

$$OPT(i, j) = \begin{cases} 0 & i \geq j - 4 \\ \max \left\{ \begin{array}{l} OPT(i, j-1) \\ \max_{\substack{i \leq t < j-4 \\ b_t, b_j \text{ complementi}}} \{ 1 + OPT(i, t-1) + OPT(t+1, j-1) \} \end{array} \right\} & \text{altrimenti} \end{cases}$$

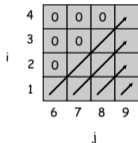
Bottom Up Dynamic Programming Over Intervals

Domanda. In che ordine risolvere i sottoproblemi?

Risposta. Dapprima gli intervalli più corti.

```
RNA( $b_1, \dots, b_n$ ) {  
  for k = 5, 6, ..., n-1  
    for i = 1, 2, ..., n-k  
      j = i + k  
      Compute M[i, j]  
  
  return M[1, n]  
}
```

k è la lunghezza
intervallo j-i

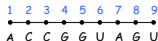


usare ricorrenza

Complessità. $O(n^3)$.

$$\text{OPT}(i, j) = \max \begin{cases} \text{OPT}(i, j-1) \\ \max_{\substack{t=i+4 \\ \Delta A \text{ complement}}} \{ 1 + \text{OPT}(i, t-1) + \text{OPT}(t+1, j-1) \} \end{cases}$$

RNA sequence ACCGGUAGU



4	0	0	0	
3	0	0		
2	0			
$i = 1$				
	$j = 6$	7	8	9

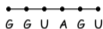
Initial values

4	0	0	0	0
3	0	0	1	
2	0	0		
$i = 1$	1			
	$j = 6$	7	8	9

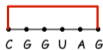
Filling in the values
for $k = 5$

$$OPT(i, j) = \max \begin{cases} OPT(i, j-1) \\ \max_{\substack{t < j-4 \\ h, h' \text{ complements}}} \{ 1 + OPT(i, t-1) + OPT(t+1, j-1) \} \end{cases}$$

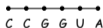
$$OPT(4, 9) = \max \begin{cases} OPT(4, 8) = 0 \\ \max_{\substack{4 \leq t < 9 \\ h, h' \text{ complements}}} \{ 1 + OPT(4, t-1) + OPT(t+1, 9) \} = 0 \end{cases}$$



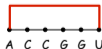
$$OPT(3, 8) = \max \begin{cases} OPT(3, 7) = 0 \\ \max_{\substack{3 \leq t < 8 \\ h, h' \text{ complements}}} \{ 1 + OPT(1, t-1) + OPT(t+1, 8) \} = 1 \end{cases}$$



$$OPT(2, 7) = \max \begin{cases} OPT(2, 6) = 0 \\ \max_{\substack{2 \leq t < 7 \\ h, h' \text{ complements}}} \{ 1 + OPT(2, t-1) + OPT(t+1, 7) \} = 0 \end{cases}$$



$$OPT(1, 6) = \max \begin{cases} OPT(1, 5) = 0 \\ \max_{\substack{1 \leq t < 6 \\ h, h' \text{ complements}}} \{ 1 + OPT(1, t-1) + OPT(t+1, 6) \} = 1 \end{cases}$$



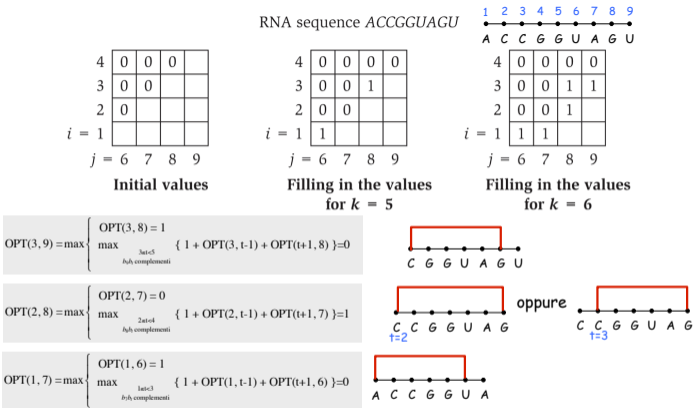


Figure 6.16 The iterations of the algorithm on a sample instance of the RNA Secondary Structure Prediction Problem.

4	0	0	0	
3	0	0		
2	0			
$i = 1$				

$j = 6 \quad 7 \quad 8 \quad 9$

Initial values

RNA sequence ACCGGUAGU

4	0	0	0	0
3	0	0	1	
2	0	0		
$i = 1$	1			

$j = 6 \quad 7 \quad 8 \quad 9$

Filling in the values
for $k = 5$

1 2 3 4 5 6 7 8 9
A C C G G U A G U

4	0	0	0	0
3	0	0	1	1
2	0	0	1	
$i = 1$	1	1		

$j = 6 \quad 7 \quad 8 \quad 9$

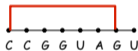
Filling in the values
for $k = 6$

4	0	0	0	0
3	0	0	1	1
2	0	0	1	1
$i = 1$	1	1	1	

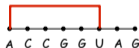
$j = 6 \quad 7 \quad 8 \quad 9$

Filling in the values
for $k = 7$

$$\text{OPT}(2, 9) = \max \begin{cases} \text{OPT}(2, 8) = 1 \\ \max_{\substack{2 \leq t < 9 \\ A_t \neq \text{complement}(A_2)}} \{ 1 + \text{OPT}(2, t-1) + \text{OPT}(t+1, 8) \} = 0 \end{cases}$$



$$\text{OPT}(1, 8) = \max \begin{cases} \text{OPT}(1, 7) = 1 \\ \max_{\substack{1 \leq t < 8 \\ A_t \neq \text{complement}(A_1)}} \{ 1 + \text{OPT}(1, t-1) + \text{OPT}(t+1, 7) \} = 1 \end{cases}$$



...

Figure 6.16 The iterations of the algorithm on a sample instance of the RNA Secondary Structure Prediction Problem.

RNA sequence ACCGGUAGU

1 2 3 4 5 6 7 8 9
A C C G G U A G U

4	0	0	0	
3	0	0		
2	0			
$i = 1$				

$j = 6 \ 7 \ 8 \ 9$

Initial values

4	0	0	0	0
3	0	0	1	
2	0	0		
$i = 1$	1			

$j = 6 \ 7 \ 8 \ 9$

Filling in the values
for $k = 5$

4	0	0	0	0
3	0	0	1	1
2	0	0	1	
$i = 1$	1	1		

$j = 6 \ 7 \ 8 \ 9$

Filling in the values
for $k = 6$

4	0	0	0	0
3	0	0	1	1
2	0	0	1	1
$i = 1$	1	1	1	

$j = 6 \ 7 \ 8 \ 9$

Filling in the values
for $k = 7$

4	0	0	0	0
3	0	0	1	1
2	0	0	1	1
$i = 1$	1	1	1	2

$j = 6 \ 7 \ 8 \ 9$

Filling in the values
for $k = 8$

$$OPT(1, 9) = \max \begin{cases} OPT(1, 8) = 1 \\ \max_{\substack{t \leq 5 \\ A_t A_{t+1} \text{ complement}}} \{ 1 + OPT(1, t-1) + OPT(t+1, 8) \} = 2 \end{cases}$$

max per $t=1$
 $OPT(1,9) = 1 + OPT(2,8)$

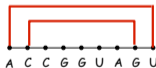


Figure 6.16 The iterations of the algorithm on a sample instance of the RNA Secondary Structure Prediction Problem.

Trovare una soluzione

RNA sequence *ACCGGUAGU*

4	0	0	0	
3	0	0		
2	0			
$i = 1$				

$j = 6 \quad 7 \quad 8 \quad 9$

Initial values

4	0	0	0	0
3	0	0	1	
2	0	0		
$i = 1$	1			

$j = 6 \quad 7 \quad 8 \quad 9$

Filling in the values
for $k = 5$

4	0	0	0	0
3	0	0	1	1
2	0	0	1	
$i = 1$	1	1		

$j = 6 \quad 7 \quad 8 \quad 9$

Filling in the values
for $k = 6$

4	0	0	0	0
3	0	0	1	1
2	0	0	1	1
$i = 1$	1	1	1	

$j = 6 \quad 7 \quad 8 \quad 9$

Filling in the values
for $k = 7$

4	0	0	0	0
3	0	0	1	1
2	0	0	1	1
$i = 1$	1	1	1	2

$j = 6 \quad 7 \quad 8 \quad 9$

Filling in the values
for $k = 8$



Figure 6.16 The iterations of the algorithm on a sample instance of the RNA Secondary Structure Prediction Problem.

Programmazione dinamica: sommario

Ricetta.

- Caratterizza struttura del problema.
- Definire ricorsivamente valore della soluzione ottimale.
- Computare valore soluzione ottimale.
- Costruire soluzione ottimale dalle informazioni computate.

Tecniche di Programmazione dinamica.

- Scelta binaria: weighted interval scheduling.
- Multi-way choice: segmented least squares.
- Aggiunta di una nuova variabile: knapsack.
- Programmazione dinamica su intervalli: RNA secondary structure.

Top-down vs. bottom-up: persone diverse hanno diverse intuizioni.

6.6 Sequence Alignment

Similarità di stringhe

Quanto sono simili due stringhe?

- **ocurrance**
- **occurrence**

o	c	u	r	r	a	n	c	e	-
o	c	c	u	r	r	e	n	c	e

6 mismatches, 1 gap

o	c	-	u	r	r	a	n	c	e
o	c	c	u	r	r	e	n	c	e

1 mismatch, 1 gap

o	c	-	u	r	r	-	a	n	c	e
o	c	c	u	r	r	e	-	n	c	e

0 mismatches, 3 gaps

Edit Distance

Applicazioni.

- Base per diff di Unix.
- Riconoscimento voce.
- Biologia Computazionale.

Edit distance. [Levenshtein 1966, Needleman-Wunsch 1970]

- Penalità per gap δ ; penalità per mismatch α_{pq} .
- Costo = somma delle penalità di gap e mismatch.

C T G A C C T A C C T

C C T G A C T A C A T

$$\alpha_{TC} + \alpha_{GT} + \alpha_{AG} + 2\alpha_{CA}$$

- C T G A C C T A C C T

C C T G A C - T A C A T

$$2\delta + \alpha_{CA}$$

Sequence Alignment

Obiettivo: Date due stringhe $X = x_1 x_2 \dots x_m$ e $Y = y_1 y_2 \dots y_n$ trovare un allineamento di costo minimo.

Definizione. Un **allineamento** M è un insieme di coppie ordinate x_i-y_j tali che ogni elemento occorre in al massimo una coppia e senza *cross*.

Definizione. La coppia x_i-y_j e $x_{i'}-y_{j'}$ **cross** se $i < i'$, ma $j > j'$.

Esempio: CTACCG vs. TACATG.

Allineamento: $M = x_2-y_1, x_3-y_2, x_4-y_3, x_5-y_4, x_6-y_6$.

x_1	x_2	x_3	x_4	x_5		x_6
C	T	A	C	C	-	G
	-	T	A	C	A	T
		y_1	y_2	y_3	y_4	y_5
						y_6

$$\text{costo}(M) = \underbrace{\sum_{(x_i, y_j) \in M} \alpha_{x_i y_j}}_{\text{mismatch}} + \underbrace{\sum_{i: x_i \text{ non accoppiato}} \delta + \sum_{j: y_j \text{ non accoppiato}} \delta}_{\text{gap}}$$

Sequence Alignment: Struttura del problema

Definizione. $OPT(i, j) = \text{min costo allineamento } x_1 x_2 \dots x_i \text{ e } y_1 y_2 \dots y_j.$

- **Caso 1:** OPT accoppia x_i - y_j .
 - mismatch per x_i - y_j + min costo allineamento $x_1 x_2 \dots x_{i-1}$ e $y_1 y_2 \dots y_{j-1}$
- **Caso 2a:** OPT lascia x_i senza accoppiamento.
 - gap per x_i e minimo costo allineamento $x_1 x_2 \dots x_{i-1}$ e $y_1 y_2 \dots y_j$
- **Caso 2b:** OPT lascia y_j senza accoppiamento.
 - gap per y_j e minimo costo allineamento $x_1 x_2 \dots x_i$ e $y_1 y_2 \dots y_{j-1}$

$$OPT(i, j) = \begin{cases} j\delta & \text{se } i=0 \\ \min \begin{cases} \alpha_{x_i y_j} + OPT(i-1, j-1) \\ \delta + OPT(i-1, j) \\ \delta + OPT(i, j-1) \end{cases} & \text{altrimenti} \\ i\delta & \text{se } j=0 \end{cases}$$

Sequence Alignment: Algorithm

```
Sequence-Alignment(m, n,  $x_1x_2\dots x_m$ ,  $y_1y_2\dots y_n$ ,  $\delta$ ,  $\alpha$ ) {  
  for i = 0 to m  
    M[0, i] =  $i\delta$   
  for j = 0 to n  
    M[j, 0] =  $j\delta$   
  
  for i = 1 to m  
    for j = 1 to n  
      M[i, j] = min( $\alpha[x_i, y_j] + M[i-1, j-1]$ ,  
                    $\delta + M[i-1, j]$ ,  
                    $\delta + M[i, j-1]$ )  
  
  return M[m, n]  
}
```

Sequence Alignment: Algorithm

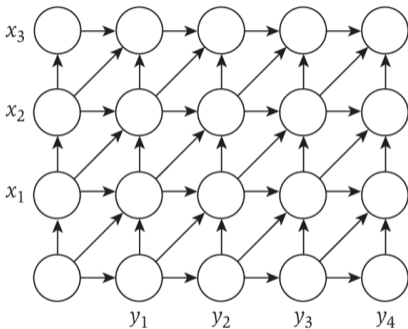


Figure 6.17 A graph-based picture of sequence alignment.

Sequence Alignment: Esempio algoritmo

Penalità per gap $\delta = 2$

Penalità per mismatch

- $\alpha_{\text{vocale,vocale}} = 1$
- $\alpha_{\text{consonante,consonante}} = 1$
- $\alpha_{\text{vocale,consonante}} = 3$

$$\begin{aligned}
 M[4,4] &= \min(3+M[3,3], 2+M[3,4], 2+M[4,3]) \\
 &= \min(8, 7, 6) \\
 &= 6
 \end{aligned}$$

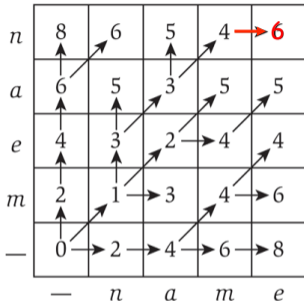


Figure 6.18 The OPT values for the problem of aligning the words *mean* to *name*.

Sequence Alignment: Esempio algoritmo

Penalità per gap $\delta = 2$

Penalità per mismatch

- $\alpha_{\text{vocale,vocale}} = 1$
- $\alpha_{\text{consonante,consonante}} = 1$
- $\alpha_{\text{vocale,consonante}} = 3$

$$\begin{aligned}
 M[4,4] &= \min(3+M[3,3], 2+M[3,4], 2+M[4,3]) \\
 &= \min(8, 7, 6) \\
 &= 6
 \end{aligned}$$

m e a n -

n - a **m** e

1 2 1 2 = 6

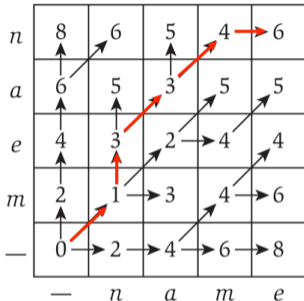


Figure 6.18 The OPT values for the problem of aligning the words *mean* to *name*.

Sequence Alignment: Complessità algoritmo

```
Sequence-Alignment(m, n, x1x2...xm, y1y2...yn, δ, α) {  
  for i = 0 to m  
    M[0, i] = iδ  
  for j = 0 to n  
    M[j, 0] = jδ  
  
  for i = 1 to m  
    for j = 1 to n  
      M[i, j] = min(α[xi, yj] + M[i-1, j-1],  
                   δ + M[i-1, j],  
                   δ + M[i, j-1])  
  
  return M[m, n]  
}
```

Analisi. $\Theta(mn)$ tempo e spazio.

Parole o frasi in italiano o inglese: $m, n \leq 10$.

Biologia Computazionale: $m = n = 100.000$.

Va bene per 10 miliardi operazioni, ma array 10GB?

Sequence Alignment: Linear Space

Domanda. Possiamo usare meno spazio?

Facile. Valore ottimo in spazio $O(m + n)$ e tempo $O(mn)$.

- Computare $OPT(i, \cdot)$ da $OPT(i-1, \cdot)$.
- Non è più facile computare un allineamento ottimale.

Teorema. [Hirschberg 1975] Allineamento ottimo in spazio $O(m + n)$ e tempo $O(mn)$.

- Combinazione di divide-and-conquer e programmazione dinamica.

6.7 Sequence Alignment in Linear Space

... ma non lo tratteremo

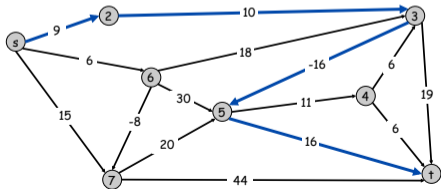
6.8 Shortest Paths

Shortest Paths

Problema del cammino più corto. Dato un grafo diretto $G = (V, E)$, con peso degli archi c_{vw} , trovare il cammino più corto dal nodo s al nodo t .

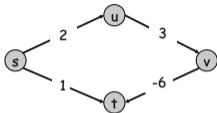
anche pesi negativi

Esempio. I nodi rappresentano agenti in ambito finanziario
 c_{vw} è il costo di una transazione nella quale si acquista da v e si vende a w .
Un cammino rappresenta una successione di transazioni.



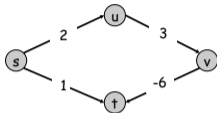
Shortest Paths: Approcci che non funzionano

Algoritmo di Dijkstra. Potrebbe non funzionare con costi negativi.

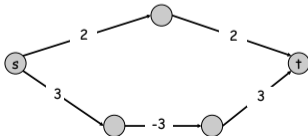


Shortest Paths: Approcci che non funzionano

Algoritmo di Dijkstra. Potrebbe non funzionare con costi negativi.

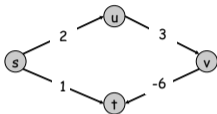


Scalare il peso degli archi. Aggiungere una costante ad ogni peso degli archi potrebbe non funzionare.

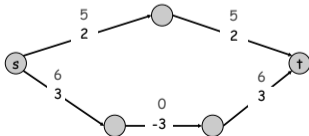


Shortest Paths: Approcci che non funzionano

Algoritmo di Dijkstra. Potrebbe non funzionare con costi negativi.

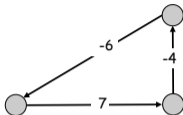


Scalare il peso degli archi. Aggiungere una costante ad ogni peso degli archi potrebbe non funzionare.



Shortest Paths: Ciclo con costo negativo

Ciclo con costo negativo.



Osservazione. Se un cammino da s a t contiene un ciclo con costo negativo, non esiste un cammino s - t più corto; comunque, ne esiste uno che è semplice.

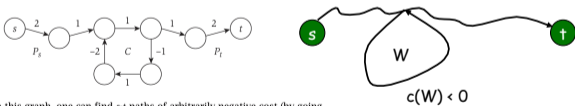


Figure 6.20 In this graph, one can find s - t paths of arbitrarily negative cost (by going around the cycle C many times).

Assumiamo che il grafo non contenga cicli con costo negativo!

Shortest Paths: Programmazione dinamica

Definizione. $OPT(i, v)$ = lunghezza del cammino P più corto $v-t$ usando al massimo i archi.

- Caso 1: P usa al massimo $i-1$ archi.
 - $OPT(i, v) = OPT(i-1, v)$
- Caso 2: P usa esattamente i archi.
 - se (v, w) è il primo arco, allora OPT usa (v, w) , e poi sceglie il miglior cammino $w-t$ usando al massimo $i-1$ archi
 - $OPT(i, v) = OPT(i-1, w) + c_{vw}$

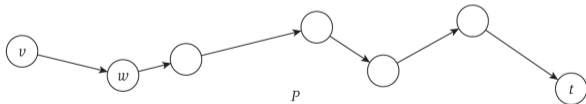


Figure 6.22 The minimum-cost path P from v to t using at most i edges.

Shortest Paths: Programmazione dinamica

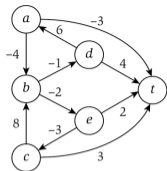
Definizione. $OPT(i, v)$ = lunghezza del cammino P più corto v-t usando al massimo i archi.

- Caso 1: P usa al massimo i-1 archi.
 - $OPT(i, v) = OPT(i-1, v)$
- Caso 2: P usa esattamente i archi.
 - se (v, w) è il primo arco, allora OPT usa (v, w) , e poi sceglie il miglior cammino w-t usando al massimo i-1 archi

$$OPT(i, v) = \begin{cases} 0 & \text{se } v = t \\ \min \left\{ OPT(i-1, v), \min_{(v,w) \in E} \{ OPT(i-1, w) + c_{vw} \} \right\} & \text{altrimenti} \\ \infty & \text{se } i = 0, v \neq t \end{cases}$$

Nota. $OPT(n-1, v)$ = lunghezza del cammino v-t più corto
(Ricorda che non ci sono cicli con costo negativo)

Shortest Paths: Esempio



(a)

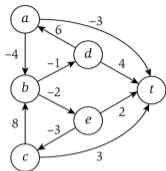
$$OPT(i, v) = \begin{cases} 0 & \text{se } v = t \\ \min \left\{ OPT(i-1, v), \min_{(v,w) \in E} \{ OPT(i-1, w) + c_{vw} \} \right\} & \text{altrimenti} \\ \infty & \text{se } i=0, v \neq t \end{cases}$$

	0	1	2	3	4	5
t	0	0	0	0	0	0
a	∞	-3	-3	-4	-6	-6
b	∞	∞	0	-2	-2	-2
c	∞	3	3	3	3	3
d	∞	4	3	3	2	0
e	∞	2	0	0	0	0

(b)

Figure 6.23 For the directed graph in (a), the Shortest-Path Algorithm constructs the dynamic programming table in (b).

Shortest Paths: Esempio



(a)

$$OPT(i, v) = \begin{cases} 0 & \text{se } v = t \\ \min \left\{ OPT(i-1, v), \min_{(v,w) \in E} \{ OPT(i-1, w) + c_{vw} \} \right\} & \text{altrimenti} \\ \infty & \text{se } i=0, v \neq t \end{cases}$$

	0	1	2	3	4	5
t	0	0	0	0	0	0
a	∞	-3	-3	-4	-6	-6
b	∞	∞	0	-2	-2	-2
c	∞	3	3	3	3	3
d	∞	4	3	3	2	0
e	∞	2	0	0	0	0

(b)

Ricorda: $OPT(n-1, v)$ = lunghezza del cammino v-t più corto

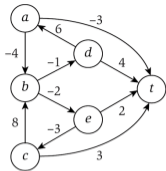
$OPT(5, a)$ = lunghezza del cammino a-t più corto

Qual è il cammino a-t più corto?

Come lo calcoliamo a partire dalla matrice?

Figure 6.23 For the directed graph in (a), the Shortest-Path Algorithm constructs the dynamic programming table in (b).

Shortest Paths: Esempio



(a)

$$OPT(i, v) = \begin{cases} 0 & \text{se } v = t \\ \min \left\{ OPT(i-1, v), \min_{(v,w) \in E} \{ OPT(i-1, w) + c_{vw} \} \right\} & \text{altrimenti} \\ \infty & \text{se } i = 0, v \neq t \end{cases}$$

	0	1	2	3	4	5
t	0	0	0	0	0	0
a	∞	-3	-3	-4	-6	-6
b	∞	∞	0	-2	-2	-2
c	∞	3	3	3	3	3
d	∞	4	3	3	2	0
e	∞	2	0	0	0	0

(b)

$OPT(n-1, v)$ = lunghezza del cammino v-t più corto

$OPT(5, a)$ = lunghezza del cammino a-t più corto

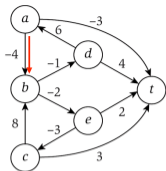
Qual è il cammino a-t più corto?

$$OPT(5, a) = \min(OPT(4, a), \min(OPT(4, t) + c_{at}, OPT(4, b) + c_{ab})) \\ = \min(-6, \min(0 - 3, -2 - 4))$$

Quindi, $OPT(5, a) = OPT(4, a)$

Figure 6.23 For the directed graph in (a), the Shortest-Path Algorithm constructs the dynamic programming table in (b).

Shortest Paths: Esempio



(a)

$$OPT(i, v) = \begin{cases} 0 & \text{se } v = t \\ \min \left\{ OPT(i-1, v), \min_{(v,w) \in E} \{ OPT(i-1, w) + c_{vw} \} \right\} & \text{altrimenti} \\ \infty & \text{se } i = 0, v \neq t \end{cases}$$

$$OPT(5, a) = \min(OPT(4, a), \min(OPT(4, t) + c_{at}, OPT(4, b) + c_{ab})) \\ = \min(-6, \min(0 - 3, -2 - 4))$$

$$OPT(4, a) = \min(OPT(3, a), \min(OPT(3, t) + c_{at}, OPT(3, b) + c_{ab})) \quad \text{a-b} \\ = \min(-4, \min(0 - 3, -2 - 4))$$

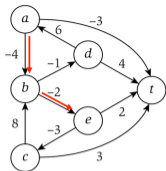
Quindi, $OPT(4, a) = OPT(3, b) + c_{ab}$

	0	1	2	3	4	5
t	0	0	0	0	0	0
a	∞	-3	-3	-4	-6	-6
b	∞	∞	0	-2	-2	-2
c	∞	3	3	3	3	3
d	∞	4	3	3	2	0
e	∞	2	0	0	0	0

(b)

Figure 6.23 For the directed graph in (a), the Shortest-Path Algorithm constructs the dynamic programming table in (b).

Shortest Paths: Esempio



(a)

$$OPT(i, v) = \begin{cases} 0 & \text{se } v = t \\ \min \left\{ OPT(i-1, v), \min_{(v,w) \in E} \{ OPT(i-1, w) + c_{vw} \} \right\} & \text{altrimenti} \\ \infty & \text{se } i=0, v \neq t \end{cases}$$

$$OPT(5, a) = \min(OPT(4, a), \min(OPT(4, t) + c_{at}, OPT(4, b) + c_{ab})) \\ = \min(-6, \min(0 - 3, -2 - 4))$$

$$OPT(4, a) = \min(OPT(3, a), \min(OPT(3, t) + c_{at}, OPT(3, b) + c_{ab})) \quad \text{a-b} \\ = \min(-4, \min(0 - 3, -2 - 4))$$

$$OPT(3, b) = \min(OPT(2, b), \min(OPT(2, e) + c_{be}, OPT(2, d) + c_{bd})) \quad \text{b-e} \\ = \min(0, \min(0 - 2, 3 - 1))$$

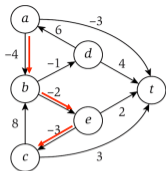
Quindi, $OPT(3, b) = OPT(2, e) + c_{be}$

	0	1	2	3	4	5
t	0	0	0	0	0	0
a	∞	-3	-3	-4	-6	-6
b	∞	∞	0	-2	-2	-2
c	∞	3	3	3	3	3
d	∞	4	3	3	2	0
e	∞	2	0	0	0	0

(b)

Figure 6.23 For the directed graph in (a), the Shortest-Path Algorithm constructs the dynamic programming table in (b).

Shortest Paths: Esempio



(a)

$$OPT(i, v) = \begin{cases} 0 & \text{se } v = t \\ \min \left\{ OPT(i-1, v), \min_{(v,w) \in E} \{ OPT(i-1, w) + c_{vw} \} \right\} & \text{altrimenti} \\ \infty & \text{se } i = 0, v \neq t \end{cases}$$

$$OPT(5, a) = \min(OPT(4, a), \min(OPT(4, t) + c_{at}, OPT(4, b) + c_{ab})) \\ = \min(-6, \min(0 - 3, -2 - 4))$$

$$OPT(4, a) = \min(OPT(3, a), \min(OPT(3, t) + c_{at}, OPT(3, b) + c_{ab})) \quad \text{a-b} \\ = \min(-4, \min(0 - 3, -2 - 4))$$

$$OPT(3, b) = \min(OPT(2, b), \min(OPT(2, e) + c_{be}, OPT(2, d) + c_{bd})) \quad \text{b-e} \\ = \min(0, \min(0 - 2, 3 - 1))$$

$$OPT(2, e) = \min(OPT(1, e), \min(OPT(1, c) + c_{ec}, OPT(1, t) + c_{et})) \quad \text{e-c} \\ = \min(2, \min(3 - 3, 0 + 2))$$

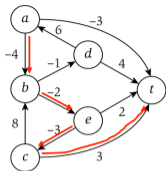
	0	1	2	3	4	5
t	0	0	0	0	0	0
a	∞	-3	-3	-4	-6	-6
b	∞	∞	0	-2	-2	-2
c	∞	3	3	3	3	3
d	∞	4	3	3	2	0
e	∞	2	0	0	0	0

(b)

Quindi, $OPT(2, e) = OPT(1, c) + c_{ec}$

Figure 6.23 For the directed graph in (a), the Shortest-Path Algorithm constructs the dynamic programming table in (b).

Shortest Paths: Esempio



(a)

	0	1	2	3	4	5
t	0	0	0	0	0	0
a	∞	-3	-3	-4	-6	-6
b	∞	∞	0	-2	-2	-2
c	∞	3	3	3	3	3
d	∞	4	3	3	2	0
e	∞	2	0	0	0	0

(b)

Figure 6.23 For the directed graph in (a), the Shortest-Path Algorithm constructs the dynamic programming table in (b).

$$OPT(i, v) = \begin{cases} 0 & \text{se } v = t \\ \min \left\{ OPT(i-1, v), \min_{(v,w) \in E} \{ OPT(i-1, w) + c_{vw} \} \right\} & \text{altrimenti} \\ \infty & \text{se } i = 0, v \neq t \end{cases}$$

$$OPT(5, a) = \min(OPT(4, a), \min(OPT(4, t) + c_{at}, OPT(4, b) + c_{ab})) \\ = \min(-6, \min(0 - 3, -2 - 4))$$

$$OPT(4, a) = \min(OPT(3, a), \min(OPT(3, t) + c_{at}, OPT(3, b) + c_{ab})) \quad \text{a-b} \\ = \min(-4, \min(0 - 3, -2 - 4))$$

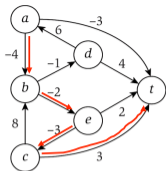
$$OPT(3, b) = \min(OPT(2, b), \min(OPT(2, e) + c_{be}, OPT(2, d) + c_{bd})) \quad \text{b-e} \\ = \min(0, \min(0 - 2, 3 - 1))$$

$$OPT(2, e) = \min(OPT(1, e), \min(OPT(1, c) + c_{ec}, OPT(1, t) + c_{et})) \quad \text{e-c} \\ = \min(2, \min(3 - 3, 0 + 2))$$

$$OPT(1, c) = \min(OPT(0, e), \min(OPT(0, b) + c_{cb}, OPT(0, t) + c_{ct})) \quad \text{c-t} \\ = \min(\infty, \min(\infty - 8, 0 + 3))$$

$$\text{Quindi, } OPT(1, c) = OPT(0, t) + c_{ct}$$

Shortest Paths: Esempio



(a)

	0	1	2	3	4	5
t	0	0	0	0	0	0
a	∞	-3	-3	-4	-6	-6
b	∞	∞	0	-2	-2	-2
c	∞	3	3	3	3	3
d	∞	4	3	3	2	0
e	∞	2	0	0	0	0

(b)

Figure 6.23 For the directed graph in (a), the Shortest-Path Algorithm constructs the dynamic programming table in (b).

$$OPT(i, v) = \begin{cases} 0 & \text{se } v = t \\ \min \left\{ OPT(i-1, v), \min_{(v,w) \in E} \{ OPT(i-1, w) + c_{vw} \} \right\} & \text{altrimenti} \\ \infty & \text{se } i = 0, v \neq t \end{cases}$$

$$OPT(5, a) = \min(OPT(4, a), \min(OPT(4, t) + c_{at}, OPT(4, b) + c_{ab})) \\ = \min(-6, \min(0 - 3, -2 - 4))$$

$$OPT(4, a) = \min(OPT(3, a), \min(OPT(3, t) + c_{at}, OPT(3, b) + c_{ab})) \quad \text{a-b} \\ = \min(-4, \min(0 - 3, -2 - 4))$$

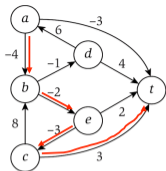
$$OPT(3, b) = \min(OPT(2, b), \min(OPT(2, e) + c_{be}, OPT(2, d) + c_{bd})) \quad \text{b-e} \\ = \min(0, \min(0 - 2, 3 - 1))$$

$$OPT(2, e) = \min(OPT(1, e), \min(OPT(1, c) + c_{ec}, OPT(1, t) + c_{et})) \quad \text{e-c} \\ = \min(2, \min(3 - 3, 0 + 2))$$

$$OPT(1, c) = \min(OPT(0, e), \min(OPT(0, b) + c_{cb}, OPT(0, t) + c_{ct})) \quad \text{c-t} \\ = \min(\infty, \min(\infty - 8, 0 + 3))$$

$$\mathbf{a \rightarrow b \rightarrow e \rightarrow c \rightarrow t} \\ \mathbf{-4 \quad -2 \quad -3 \quad +3 \quad = -6}$$

Shortest Paths: Esempio



(a)

$$OPT(i, v) = \begin{cases} 0 & \text{se } v = t \\ \min \left\{ OPT(i-1, v), \min_{(v,w) \in E} \{ OPT(i-1, w) + c_{vw} \} \right\} & \text{altrimenti} \\ \infty & \text{se } i = 0, v \neq t \end{cases}$$

$$OPT(5, a) = \min(OPT(4, a), \min(OPT(4, t) + c_{at}, OPT(4, b) + c_{ab})) \\ = \min(-6, \min(0 - 3, -2 - 4))$$

$$OPT(4, a) = \min(OPT(3, a), \min(OPT(3, t) + c_{at}, OPT(3, b) + c_{ab})) \quad \text{a-b} \\ = \min(-4, \min(0 - 3, -2 - 4))$$

$$OPT(3, b) = \min(OPT(2, b), \min(OPT(2, e) + c_{be}, OPT(2, d) + c_{bd})) \quad \text{b-e} \\ = \min(0, \min(0 - 2, 3 - 1))$$

$$OPT(2, e) = \min(OPT(1, e), \min(OPT(1, c) + c_{ec}, OPT(1, t) + c_{et})) \quad \text{e-c} \\ = \min(2, \min(3 - 3, 0 + 2))$$

$$OPT(1, c) = \min(OPT(0, e), \min(OPT(0, b) + c_{cb}, OPT(0, t) + c_{ct})) \quad \text{c-t} \\ = \min(\infty, \min(\infty - 8, 0 + 3))$$

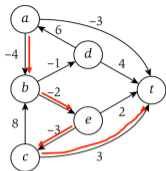
	0	1	2	3	4	5
t	0	0	0	0	0	0
a	∞	-3	-3	-4	-6	-6
b	∞	∞	0	-2	-2	-2
c	∞	3	3	3	3	3
d	∞	4	3	3	2	0
e	∞	2	0	0	0	0

(b)

Figure 6.23 For the directed graph in (a), the Shortest-Path Algorithm constructs the dynamic programming table in (b).

Shortest Paths: Esempio

Trovare il cammino più corto.
Qual'è l'informazione che possiamo ricordare?



(a)

	0	1	2	3	4	5
t	0	0	0	0	0	0
a	∞	-3	-3	-4	-6	-6
b	∞	∞	0	-2	-2	-2
c	∞	3	3	3	3	3
d	∞	4	3	3	2	0
e	∞	2	0	0	0	0

(b)

Figure 6.23 For the directed graph in (a), the Shortest-Path Algorithm constructs the dynamic programming table in (b).

$$\text{OPT}(5,a) = \min(\text{OPT}(4,a), \min(\text{OPT}(4,t)+c_{at}, \text{OPT}(4,b)+c_{ab})) \\ = \min(-6, \min(0-3, -2-4))$$

$$\text{OPT}(4,a) = \min(\text{OPT}(3,a), \min(\text{OPT}(3,t)+c_{at}, \text{OPT}(3,b)+c_{ab})) \quad \mathbf{a-b} \\ = \min(-4, \min(0-3, -2-4))$$

$$\text{OPT}(3,b) = \min(\text{OPT}(2,b), \min(\text{OPT}(2,e)+c_{be}, \text{OPT}(2,d)+c_{bd})) \quad \mathbf{b-e} \\ = \min(0, \min(0-2, 3-1))$$

$$\text{OPT}(2,e) = \min(\text{OPT}(1,e), \min(\text{OPT}(1,c)+c_{ec}, \text{OPT}(1,t)+c_{et})) \quad \mathbf{e-c} \\ = \min(2, \min(3-3, 0+2))$$

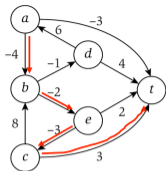
$$\text{OPT}(1,c) = \min(\text{OPT}(0,e), \min(\text{OPT}(0,b)+c_{cb}, \text{OPT}(0,t)+c_{ct})) \quad \mathbf{c-t} \\ = \min(\infty, \min(\infty-8, 0+3))$$

Shortest Paths: Esempio

Trovare il cammino più corto.

Qual'è l'informazione che possiamo ricordare?

Mantenere un "successore" per ogni valore della tabella.



(a)

	0	1	2	3	4	5
t	0	0	0	0	0	0
a	∞	-3	-3	-4	-6	-6
b	∞	∞	0	-2	-2	-2
c	∞	3	3	3	3	3
d	∞	4	3	3	2	0
e	∞	2	0	0	0	0

(b)

$$\begin{aligned} \text{OPT}(5,a) &= \min(\text{OPT}(4,a), \min(\text{OPT}(4,t)+c_{at}, \text{OPT}(4,b)+c_{ab})) \\ &= \min(-6, \min(0-3, -2-4)) \end{aligned}$$

$$\begin{aligned} \text{OPT}(4,a) &= \min(\text{OPT}(3,a), \min(\text{OPT}(3,t)+c_{at}, \text{OPT}(3,b)+c_{ab})) & \mathbf{a-b} \\ &= \min(-4, \min(0-3, -2-4)) \end{aligned}$$

$$\begin{aligned} \text{OPT}(3,b) &= \min(\text{OPT}(2,b), \min(\text{OPT}(2,e)+c_{be}, \text{OPT}(2,d)+c_{bd})) & \mathbf{b-e} \\ &= \min(0, \min(0-2, 3-1)) \end{aligned}$$

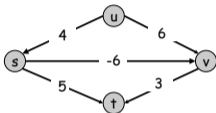
$$\begin{aligned} \text{OPT}(2,e) &= \min(\text{OPT}(1,e), \min(\text{OPT}(1,c)+c_{ec}, \text{OPT}(1,t)+c_{et})) & \mathbf{e-c} \\ &= \min(2, \min(3-3, 0+2)) \end{aligned}$$

$$\begin{aligned} \text{OPT}(1,c) &= \min(\text{OPT}(0,e), \min(\text{OPT}(0,b)+c_{cb}, \text{OPT}(0,t)+c_{ct})) & \mathbf{c-t} \\ &= \min(\infty, \min(\infty-8, 0+3)) \end{aligned}$$

Figure 6.23 For the directed graph in (a), the Shortest-Path Algorithm constructs the dynamic programming table in (b).

Shortest Paths: Esercizio

Si esegua l'algoritmo di programmazione dinamica Shortest-Path(G, t) per il calcolo dei cammini minimi sul grafo G



Si chiariscano i passi effettuati evidenziando i valori della matrice OPT costruita dall'algoritmo

	0	1	2	3
t	0	0	0	0
u	∞			
v	∞			
s	∞			

Si descriva come ottenere il cammino di costo minimo dal nodo u facendo uso della matrice OPT e chiarendo i passi effettuati.

Shortest Paths: Implementazione

```
Shortest-Path(G, t) {  
  array M[0...n-1,V]  
  foreach node v ∈ V  
    M[0, v] ← ∞  
  M[0, t] ← 0  
  
  for i = 1 to n-1  
    foreach node v ∈ V  
      M[i, v] ← M[i-1, v]  
    foreach edge (v, w) ∈ E  
      M[i, v] ← min { M[i, v], M[i-1, w] + cvw }  
}
```

Shortest Paths: Implementazione

```
Shortest-Path(G, t) {  
  array M[0..n-1, V]  
  foreach node v ∈ V  
    M[0, v] ← ∞  
  M[0, t] ← 0  
  
  for i = 1 to n-1  
    foreach node v ∈ V  
      M[i, v] ← M[i-1, v]  
    foreach edge (v, w) ∈ E  
      M[i, v] ← min { M[i, v], M[i-1, w] + cvw }  
}
```

Analisi. Tempo $\Theta(mn)$, spazio $\Theta(n^2)$.

Shortest Paths: Miglioramenti pratici

Miglioramenti pratici.

- Mantenere solo un array $M[v]$ = cammino più corto $v-t$ che abbiamo trovato finora. ("i" è solo un contatore)

$$M(i, v) = \min \left\{ M(i-1, v), \min_{(v,w) \in E} \{ M(i-1, w) + c_{vw} \} \right\}$$

- la relazione "si semplifica" :

$$M(v) = \min \left\{ M(v), \min_{(v,w) \in E} \{ M(w) + c_{vw} \} \right\}$$

```
Shortest-Path(G, t) {  
  array M[V]  
  foreach node v ∈ V  
    M[v] ← ∞  
  M[t] ← 0  
  
  for i = 1 to n-1  
    foreach edge (v, w) ∈ E  
      M[v] ← min { M[v], M[w] + cvw }  
}
```

Shortest Paths: Miglioramenti pratici

Miglioramenti pratici.

- Mantenere solo un array $M[v]$ = cammino più corto $v-t$ che abbiamo trovato finora.
- Nessuna necessità di controllare archi della forma (v, w) a meno che $M[w]$ è cambiato nell'iterazione precedente.

Teorema. Durante l'algoritmo, $M[v]$ è la lunghezza di un cammino $v-t$, e dopo i round di aggiornamenti, il valore $M[v]$ non è maggiore della lunghezza del cammino $v-t$ più corto usando $\leq i$ archi.

Impatto totale.

- Memoria: $O(m + n)$.
- Complessità tempo: caso peggiore $O(mn)$, ma veloce in pratica.

Esercizio: Longest Common Subsequence

Problema: Date 2 sequenze, $X = x_1 \dots x_m$ e $Y = y_1 \dots y_n$, trovare una sottosequenza comune la cui lunghezza è massima.

springtime



printing

ncaa tournament



north carolina

basketball



krzyzewski

Algoritmo naïve: $\Theta(n2^m)$

(per ogni sottosequenza di X , verificare se è anche sottosequenza di Y)

Esercizio: Longest Common Subsequence

Problema: Date 2 sequenze, $X = x_1 \dots x_m$ e $Y = y_1 \dots y_n$, trovare una sottosequenza comune la cui lunghezza è massima.

springtime



printing

ncaa tournament



north carolina

basketball



krzyzewski

Algoritmo naïve: $\Theta(n2^m)$

(per ogni sottosequenza di X , verificare se è anche sottosequenza di Y)

Descrivere ed analizzare un algoritmo di programmazione dinamica

Suggerimento: la tecnica è simile a quella per Sequence Alignment

Sia $Z = z_1 \dots z_k$ una LCS di X e Y

1. Se $x_m = y_n$, allora $z_k = x_m = y_n$ e Z_{k-1} è una LCS di X_{m-1} e Y_{n-1}

2. Se $x_m \neq y_n$, allora

- $z_k \neq x_m$ e Z è una LCS di X_{m-1} e Y oppure
- $z_k \neq y_n$ e Z è una LCS di X e Y_{n-1}

Notazione:

$X_i = x_1 \dots x_i$ è un prefisso di X

$Y_i = y_1 \dots y_i$ è un prefisso di Y

Longest Common Subsequence: soluzione ricorsiva

Sia $Z = z_1 \dots z_k$ una LCS di X e Y

1. Se $x_m = y_n$, allora $z_k = x_m = y_n$ e Z_{k-1} è una LCS di X_{m-1} e Y_{n-1}
2. Se $x_m \neq y_n$, allora
 - $z_k \neq x_m$ e Z è una LCS di X_{m-1} e Y oppure
 - $z_k \neq y_n$ e Z è una LCS di X e Y_{n-1}

Definiamo $c[i,j]$ = lunghezza della LCS di X_i e Y_j

La soluzione del problema è data da $c[m,n]$.

$$c[i,j] = \begin{cases} 0 & \text{se } i = 0 \text{ oppure } j = 0, \\ c[i-1, j-1] + 1 & \text{se } i, j > 0 \text{ e } x_i = y_j, \\ \max(c[i-1, j], c[i, j-1]) & \text{se } i, j > 0 \text{ e } x_i \neq y_j. \end{cases}$$

Longest Common Subsequence: algoritmo iterativo

```
LCS (X, Y)
m ← length[X]
n ← length[Y]
for i ← 1 to m
  do c[i,0] ← 0
for j ← 0 to n
  do c[0,j] ← 0
for i ← 1 to m
  do for j ← 1 to n
    do if  $x_i = y_j$ 
      then c[i,j] ← c[i-1,j-1] + 1
         b[i,j] ← “\”
    else if c[i-1,j] ≥ c[i,j-1]
      then c[i,j] ← c[i-1,j]
         b[i,j] ← “↑”
    else c[i,j] ← c[i,j-1]
         b[i,j] ← “←”
return c and b
```

$b[i,j]$ è un puntatore al sottoproblema usato per risolvere la LCS di X_i e Y_j

I valori $b[i,j]$ sono utili per calcolare la sottosequenza di lunghezza massima

Il valore $c[m,n]$ è la lunghezza della LCS di X e Y.

Complessità $O(mn)$

Longest Common Subsequence: esempio

		j	0	1	2	3	4	5	6
i	y_j		B	D	C	A	B	A	
	x_i	0	0	0	0	0	0	0	0
0	A	0	0	0	0	1	1	1	
1	B	0	1	1	1	1	2	2	
2	C	0	1	1	2	2	2	2	
3	B	0	1	1	2	2	3	3	
4	D	0	1	2	2	2	3	3	
5	A	0	1	2	2	3	3	4	
6	B	0	1	2	2	3	4	4	

Figure 15.6 The c and b tables computed by LCS-LENGTH on the sequences $X = \langle A, B, C, B, D, A, B \rangle$ and $Y = \langle B, D, C, A, B, A \rangle$. The square in row i and column j contains the value of $c[i, j]$ and the appropriate arrow for the value of $b[i, j]$. The entry 4 in $c[7, 6]$ —the lower right-hand corner of the table—is the length of an LCS $\langle B, C, B, A \rangle$ of X and Y . For $i, j > 0$, entry $c[i, j]$ depends only on whether $x_i = y_j$ and the values in entries $c[i - 1, j]$, $c[i, j - 1]$, and $c[i - 1, j - 1]$, which are computed before $c[i, j]$. To reconstruct the elements of an LCS, follow the $b[i, j]$ arrows from the lower right-hand corner; the path is shaded. Each “↖” on the path corresponds to an entry (highlighted) for which $x_i = y_j$ is a member of an LCS.

Longest Common Subsequence: costruzione di una LCS

```
PRINT-LCS ( $b, X, i, j$ )
if  $i = 0$  or  $j = 0$ 
    then return
if  $b[i, j] = "\backslash"$ 
    then PRINT-LCS( $b, X, i-1, j-1$ )
    print  $x_i$ 
    else if  $b[i, j] = "\uparrow"$ 
        then PRINT-LCS( $b, X, i-1, j$ )
else PRINT-LCS( $b, X, i, j-1$ )
```

- Chiamata iniziale: PRINT-LCS (b, X, i, j)
- Quando $b[i, j] = "\backslash"$, la LCS viene estesa di un carattere.
- Running Time: $O(m \times n)$

Riepilogo Cap. 6, Dynamic Programming

- 6.1 Weighted Interval Scheduling
- 6.2 Principles of Dynamic Programming
- 6.3 Segmented Least Squares
- 6.4 Knapsack Problem
- 6.5 RNA Secondary Structure
- 6.6 Sequence Alignment
- 6.7 Sequence Alignment in Linear Space (no)
- 6.8 Shortest Paths in a Graph
- 6.9 Shortest Paths and Distance Vector Protocols (no)
- 6.10 Negative Cycles in a Graph (no)

Esercizio:

- Longest Common Subsequence