

Lezione 17 06/12/2023

Occorrenza approssimata (recap)

DEFINIZIONE

Una posizione i del testo T tale che esista una sottostringa $S=T[i-L+1,i]$ tale che $ED(P,S) \leq k$, è un'occorrenza approssimata di P in T

ESEMPIO:

$T = \text{bbac}\downarrow\text{cbb}\text{aac}$, $k = 1$ $S_1 = \text{cbb}$
 $P = \text{cbb}$ $ED(P, S_1) = 0 \leq k$
 $i = 7 \rightarrow$ occorrenza approssimata!

Testo T , pattern P , soglia di errore 1.
Dove c'è la freccia è un'occorrenza approssimata? sì.

ESEMPIO:

$T = \text{bbac}\downarrow\text{cbb}\text{aac}$, $k = 1$ $S_2 = \text{cbb}$
 $P = \text{c}\text{~~z~~bb}$ $ED(P, S_2) = 1 \leq k$
 $i = 7 \rightarrow$ occorrenza approssimata!

Anche qui è un'occorrenza approssimata perché abbiamo $k=1$

String matching approssimato

INPUT:

Un **testo** T di lunghezza n , un **pattern** P di lunghezza m , definiti su alfabeto Σ e una **soglia** k di errore

OUTPUT:

tutte le posizioni i di T in cui finisce una sottostringa S tale che $ED(P, S) \leq k$

Ricerca esatta con WM

1. Preprocessing del pattern P in tempo $\theta(|\Sigma| + m)$
 \rightarrow calcolo di $|\Sigma|$ words di m bit
2. Scansione del testo T in tempo $\theta(kn)$ per cercare le occorrenze approssimate di P

Nella fase di preprocessing si calcolava una tabella di parole di m bit, una per ogni simbolo dell'alfabeto, che dicevano come si distribuiscono i simboli. Quindi 1) è uguale a prima.

Nella scansione del testo però ora abbiamo anche k (soglia di errore) quindi non è una sola scansione, sono più scansioni successive e quindi la complessità è più alta.

Definizione di $word D_i^h$
 $(0 \leq h \leq k, 0 \leq i \leq n)$

$P[1,j] = \text{suff}_h(T[1,i])$

significa

" $P[1,j]$ è uguale a un suffisso di $T[1,i]$ a meno di h errori"

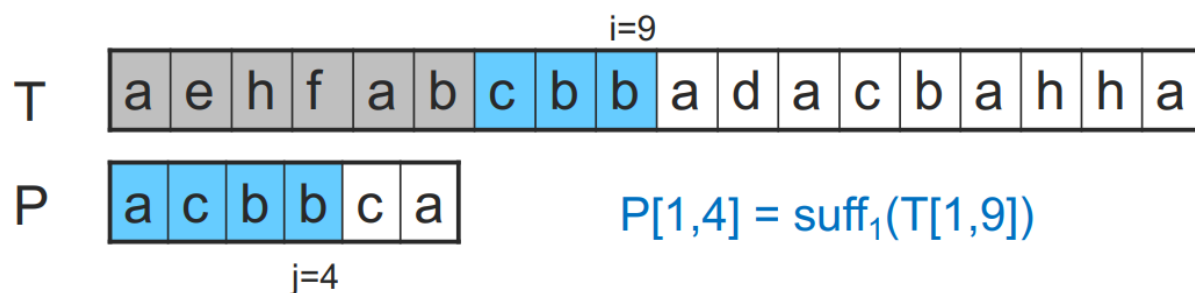
Cioè, riesco a trovare un suffisso S di $T[1,i]$ tale che:

$$ED(S, P[1,j]) \leq h$$

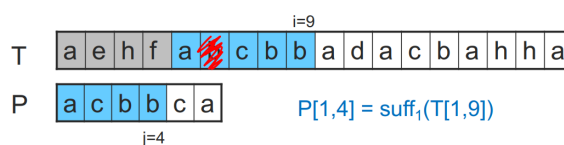
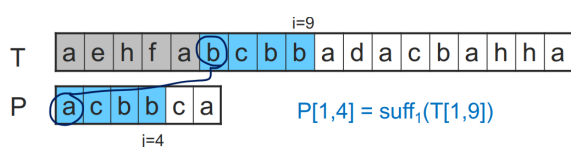
Dobbiamo definire delle parole di m bit ma più complicate. Avremo una parola, oltre alle posizioni del testo, per ogni possibile valore di k . Quindi i varia tra 0 e n , mentre h varia tra errore 0 e k . Avremo quindi $(k+1)(n+1)$ parole.

Quando scriviamo la prima riga blu, significa che possiamo trovare una sottostringa del testo che finisce in i che ha distanza di edit con il prefisso del pattern lungo j che è minore o uguale ad h .

Nell'ultima riga suffisso = sottostringa che finisce in i .



Troviamo una sottostringa che finisce in 9, che ha distanza di edit con il prefisso del pattern (blu del pattern) che è minore uguale a 1. cbb ha distanza di edit con $acbb$ 1.



sono corrette anche queste due stringhe.

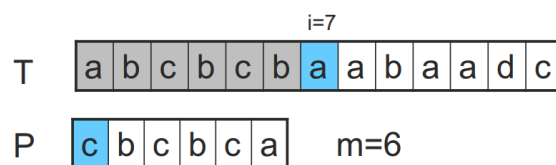
Tornando alla definizione di word D:

Dati P lungo m e T lungo n e una soglia k di errore, D^h_i ($0 \leq h \leq k, 0 \leq i \leq n$) è una *word* di m bit tale che:

$$D^h_i[j] = 1 \Leftrightarrow P[1,j] = \text{suff}_h(T[1,i])$$

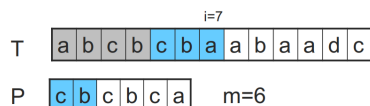
Esempio

$$D^1_7 = \underline{1}11111$$

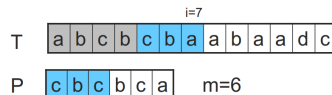


Il primo bit è 1 perchè il prefisso lungo 1 del pattern occorre come suffisso del prefisso lungo 7 nel testo con al più un errore. (Il prefisso lungo 1 del pattern ha distanza di edit con un sottostringa del testo che finisce in 7 che è minore di 1)
Il primo bit 1 è perché la distanza di edit è 1.

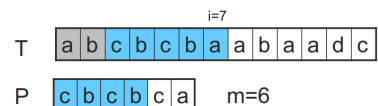
$$D^1_7 = \underline{1}11111$$



$$D^1_7 = 11111\underline{1}$$



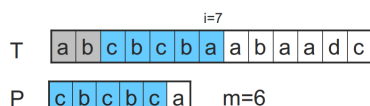
$$D^1_7 = 1111\underline{1}1$$



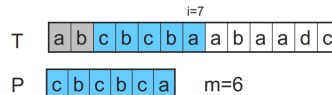
Anche qui abbiamo distanza di edit 1.

Cancelliamo la a (errore), siamo nella soglia.

$$D^1_7 = 11111\underline{1}$$



$$D^1_7 = 11111\underline{1}$$



L'errore può essere sia sul patter che sul testo, può essere fatto con una cancellazione, un inserimento, una sostituzione.

Tornando alla definizione di word D

Dati P lungo m e T lungo n e una soglia k di errore, D_i^h ($0 \leq h \leq k$, $0 \leq i \leq n$) è una word di m bit tale che:

$$D_i^h[j] = 1 \Leftrightarrow P[1,j] = \text{suff}_h(T[1,i])$$

D_i^0 ($h=0$) è per definizione la word D_i di BYG.

D_i^0 ($i=0$) è tale che:

$$j \leq h \Rightarrow D_i^0[j] = 1 \quad (P[1,j] = \text{suff}_h(T[1,0]))$$

$$j > h \Rightarrow D_i^0[j] = 0 \quad (P[1,j] \text{ non è } \text{suff}_h(T[1,0]))$$

$$D_i^0 = \underbrace{1 \ 1 \ \dots 1}_{h \text{ bit a } 1} \underbrace{0 \ \dots 0}_{(m-h) \text{ bit a } 0}$$

$$D_i^h[m] = 1 \Leftrightarrow P[1,m] = \text{suff}_h(T[1,i])$$

In particolare:

$$D_i^0[m] = 1 \Leftrightarrow P[1,m] = \text{suff}_0(T[1,i])$$

$$D_i^k[m] = 1 \Leftrightarrow P[1,m] = \text{suff}_k(T[1,i])$$

D a livello 0 abbiamo quindi tutte le parole esatte, senza errore. Quindi è lo stesso algoritmo BYG.

Quando $i=0$ il bit j sarà 1 se $j \leq h$. Anche se la porzione di testo analizzata è nulla, il bit sarà comunque 1 quando la lunghezza del prefisso del pattern è \leq ad h, ovvero gli errori massimi accettati.

Se invece $j > h$ abbiamo solo bit a 0.

Questo è il segnale per l'occorrenza approssimata.

$$D_i^h[m] = 1 \Leftrightarrow i \text{ è occorrenza approssimata per soglia } h$$

In particolare:

$$D_i^0[m] = 1 \Leftrightarrow i-m+1 \text{ è occorrenza esatta}$$

$$D_i^k[m] = 1 \Leftrightarrow i \text{ è occorrenza approssimata}$$

Per un h generico che varia tra 0 e k questa è l'occorrenza approssimata.

$$D_i^h[m] = 1 \Rightarrow D_i^{h+1}[m] = 1 \Rightarrow D_i^{h+2}[m] = 1 \Rightarrow \dots \Rightarrow D_i^k[m] = 1$$

Se il bit m della parola dove h è 1, come sarà lo stesso bit della parola con h successiva? Sarà sempre 1 perché sto semplicemente concedendo un errore in più. E così via.

$$D_i^h[m] = 1 \quad \Rightarrow \quad D_i^k[m] = 1$$

Quindi se abbiamo una parola per un h intermedio che ha l'ultimo bit a 1, allora tutte quelle con h maggiore avranno l'ultimo bit a 1.

Quindi le occorrenze approssimate saranno trovate guardando tutte le parole al livello k.

Scansione del testo

Il testo T viene scandito per $k+1$ volte:

- **Iterazione 0** → calcolo delle *words* D^0_i
- **Iterazione 1** → calcolo delle *words* D^1_i
- **Iterazione 2** → calcolo delle *words* D^2_i
- ...
- **Iterazione h** → calcolo delle *words* D^h_i
- ...
- **Iterazione k** → calcolo delle *words* D^k_i
- **All'Iterazione k**, se $D^k_i[m] = 1$, viene prodotta in output l'occorrenza approssimata i

Ad ogni iterazione cambia la h .

Tutte le occorrenze trovate nei livelli precedenti saranno trovate all'ultimo livello.

Iterazione $h=0$: calcolo della *word* D^0_i per i da 0 a n
→ algoritmo BYG

Word D^h_0 per h da 1 a k ?
 D^h_0 ha primi h bit a 1 e i restanti $(m-h)$ bit a 0

QUESTIONE DA RISOLVERE:

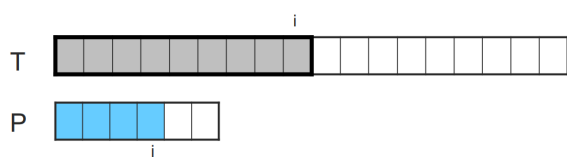
calcolare la parola D^h_i per $i > 0$ e $h > 0$

Come calcoliamo la parola generica per $i > 0$ e $h > 0$?

[Calcolo di $D^h_i[j]$ ($h > 0, i > 0$)

$j > 1$

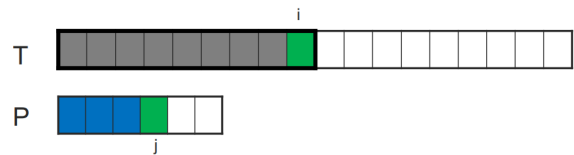
$$D^h_i[j] = 1 \Leftrightarrow P[1,j] = \text{suff}_h(T[1,i])$$



C'è una condizione che si verifica che mi rende vera l'affermazione in blu, e di conseguenza il bit 1?

$$P[1,j-1] = \text{suff}_h(T[1,i-1]) \text{ AND } P[j] = T[i]$$

se il prefisso accorciato del pattern ha distanza di edit \leq ad h e gli ultimi due caratteri devono essere uguali.



$$D^h_i[j] = 1 \Leftarrow P[1,j-1] = \text{suff}_h(T[1,i-1]) \text{ AND } P[j] = T[i]$$

Non è un sse perchè se l'AND è falso, il bit potrebbe essere comunque a 1, ci arriviamo.

Possiamo sostituire il pezzo blu con la D di i-1 e il suo bit j-1.

$$D^h_i[j] = 1 \Leftarrow D^{h-1}_{i-1}[j-1] = 1 \text{ AND } B_{T[i]}[j] = 1$$

Sostituiamo il secondo pezzo con il bit j-esimo della parola B (di preprocessing)

$$D^h_i[j] \Leftarrow D^{h-1}_{i-1}[j-1] \text{ AND } B_{T[i]}[j]$$

CASO1

Arriviamo quindi a questa affermazione finale, che è il caso 1.

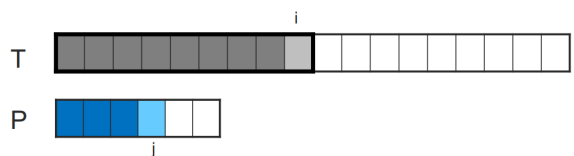
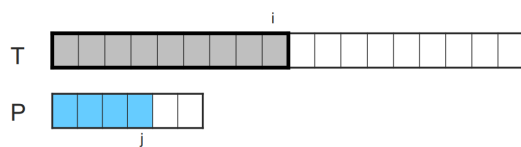
[Calcolo di $D^h_i[j]$ ($h > 0, i > 0$)

j > 1

$$D^h_i[j] = 1 \Leftrightarrow P[1,j] = \text{suff}_h(T[1,i]) \Leftrightarrow ?$$

Il caso 2 richiede un errore in meno. Quindi l'ultimo bit può anche essere diverso, non lo consideriamo.

$$D^h_i[j] = 1 \Leftarrow P[1,j-1] = \text{suff}_{h-1}(T[1,i-1])$$



$$D^h_i[j] \Leftarrow D^{h-1}_{i-1}[j-1]$$

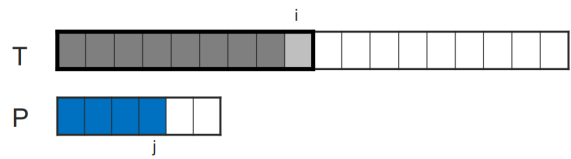
CASO2

Troviamo quindi il secondo caso. C'è anche un terzo caso:

[Calcolo di $D^h_i[j]$ ($h > 0, i > 0$)

$j > 1$

$$D^h_i[j] = 1 \Leftrightarrow P[1,j] = \text{suff}_h(T[1,i]) \Leftrightarrow P[1,j] = \text{suff}_{h-1}(T[1,i-1])$$



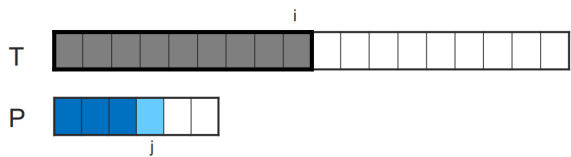
$$D^h_i[j] \Leftarrow D^{h-1}_{i-1}[j]$$

CASO3

[Calcolo di $D^h_i[j]$ ($h > 0, i > 0$)

$j > 1$

$$D^h_i[j] = 1 \Leftrightarrow P[1,j] = \text{suff}_h(T[1,i]) \Leftrightarrow P[1,j-1] = \text{suff}_{h-1}(T[1,i])$$



$$D^h_i[j] \Leftarrow D^{h-1}_i[j-1]$$

CASO4

Riassumendo:

$$1 \quad D^h_i[j] \Leftarrow (D^{h-1}_{i-1}[j-1] \text{ AND } B_{T[i]}[j])$$

OR

$$2 \quad D^h_i[j] \Leftarrow D^{h-1}_{i-1}[j-1]$$

OR

$$3 \quad D^h_i[j] \Leftarrow D^{h-1}_{i-1}[j]$$

OR

$$4 \quad D^h_i[j] \Leftarrow D^{h-1}_i[j-1]$$

Il bit j-esimo della parola h e i, sarà 1 quando almeno 1 di questi 4 casi è 1.

$$D^h_i[j] =$$

$$(D^{h-1}_{i-1}[j-1] \text{ AND } B_{T[i]}[j])$$

OR

$$D^{h-1}_{i-1}[j-1]$$

OR

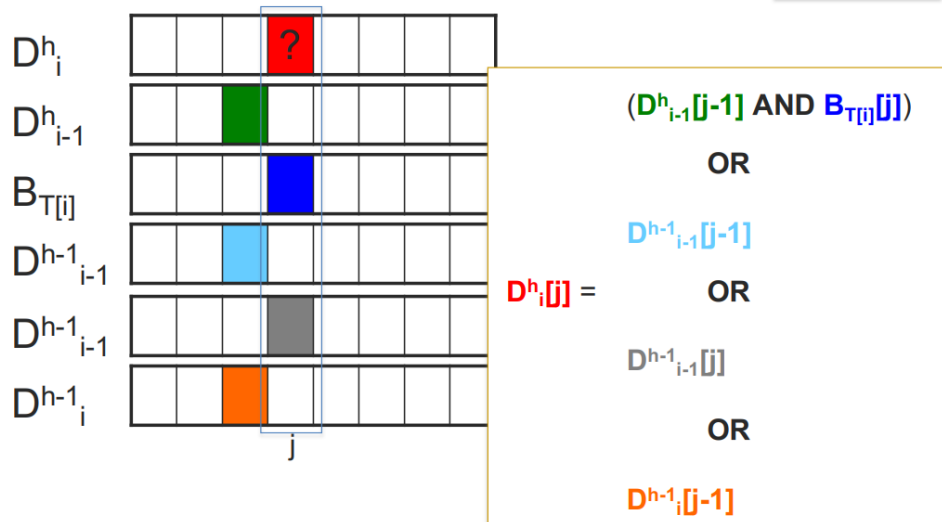
$$D^{h-1}_{i-1}[j]$$

OR

$$D^{h-1}_i[j-1]$$

[Calcolo di $D_i^h[j]$ ($h > 0, i > 0$)

$j > 1$



Passiamo al calcolo del primo bit della parola generica ($j=1$)

[Calcolo di $D_i^h[1]$ ($h > 0, i > 0$)

$j = 1$

$D_i^h[1] = 1$ per $h > 0, i > 0$

Se $h > 0$ sarà sempre 1 perché gli errori sono concessi.

Quando invece siamo nella posizione j

$$\begin{aligned}
 D_i^h[j] &= (D_{i-1}^h[j-1] \text{ AND } B_{T[i]}[j]) \text{ OR } D_{i-1}^{h-1}[j-1] \text{ OR } D_{i-1}^{h-1}[j] \text{ OR } D_i^{h-1}[j-1] \\
 D_i^h[1] &= (D_{i-1}^h[0] \text{ AND } B_{T[i]}[1]) \text{ OR } D_{i-1}^{h-1}[0] \text{ OR } D_{i-1}^{h-1}[1] \text{ OR } D_i^{h-1}[0]
 \end{aligned}$$

L'ultimo caso (posizione 0) non esiste. Quindi sostituisco brutalmente quella posizione 0 con 1.

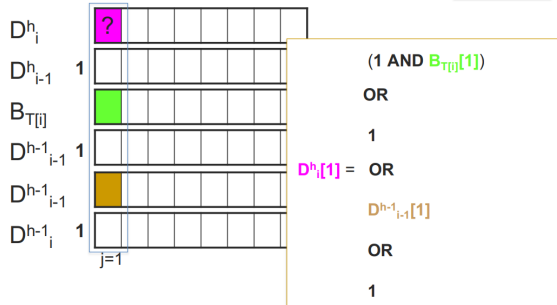
$D^h_i[1] = 1$ per $h > 0, i > 0$

$$D^h_i[1] = \begin{matrix} (1 \text{ AND } B_{T[i]}[1]) \\ \text{OR} \\ 1 \\ \text{OR} \\ D^{h-1}_{i-1}[1] \\ \text{OR} \\ 1 \end{matrix}$$

Sostituisco 1 al posto di $D^{h-1}_{i-1}[0]$, $D^{h-1}_{i-1}[0]$ e $D^{h-1}_i[0]$

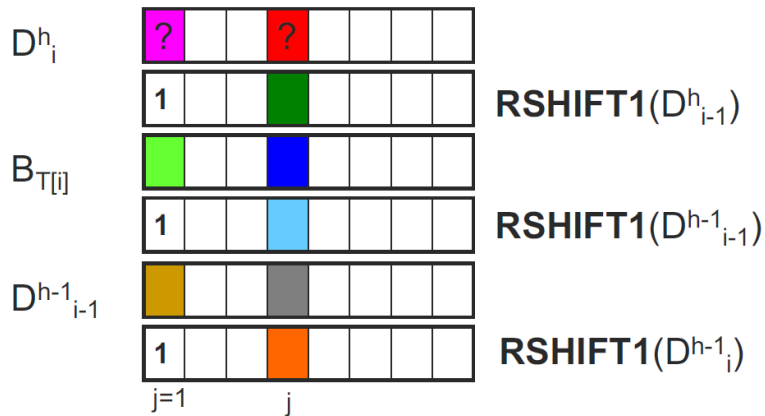
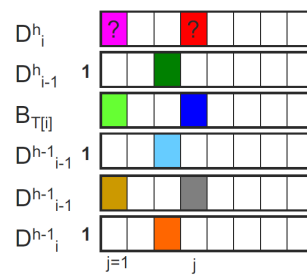
Calcolo di $D^h_i[1]$ ($h > 0, i > 0$)

$j = 1$



Calcolo di $D^h_i[j]$ ($h > 0, i > 0$)

$j = 1$
 $j > 1$



Quindi per $j > 1$

E per $j = 1$

$$\begin{array}{lcl}
 & (\text{RSHIFT1}(D_{i-1}^h)[j] \text{ AND } B_{T[i]}[j]) & (\text{RSHIFT1}(D_{i-1}^h)[1] \text{ AND } B_{T[i]}[1]) \\
 & \text{OR} & \text{OR} \\
 & \text{RSHIFT1}(D_{i-1}^{h-1})[j] & \text{RSHIFT1}(D_{i-1}^{h-1})[1] \\
 D_i^h[j] = & \text{OR} & D_i^h[1] = \\
 & D_{i-1}^{h-1}[j] & D_{i-1}^{h-1}[1] \\
 & \text{OR} & \text{OR} \\
 & \text{RSHIFT1}(D_{i-1}^{h-1})[j] & \text{RSHIFT1}(D_{i-1}^{h-1})[1]
 \end{array}$$

Quindi li possiamo mettere insieme ($j \geq 1$):

$$\begin{array}{lcl}
 & (\text{RSHIFT1}(D_{i-1}^h)[j] \text{ AND } B_{T[i]}[j]) & \\
 & \text{OR} & \\
 & \text{RSHIFT1}(D_{i-1}^{h-1})[j] & \\
 D_i^h[j] = & \text{OR} & \\
 & D_{i-1}^{h-1}[j] & \\
 & \text{OR} & \\
 & \text{RSHIFT1}(D_{i-1}^{h-1})[j] &
 \end{array}$$

1. Inizializzazione della parola $D_0^0 = 00 \dots 0$
2. Calcolo di D_i^0 per i da 1 a n (algoritmo BYG)
1. Per h da 1 a k :
 1. Inizializzazione di D_0^h
 2. Per i da 1 a n

Calcolo di D_i^h

Se $h=k$ AND $D_i^k[m] = 1$

output i

Questo è quindi l'algoritmo per la scansione del testo.

Non chiederà mai la formula di calcolo per la parola, non chiede di fare un'esecuzione su un testo. All'esame può fare domande precise sulla parola singola, ora facciamo esercizi, ci interessa l'idea dell'algoritmo, come sono fatte le parole D con i e h .

Esercizi

Esercizio 1

La parola 0100 è la parola D^0_i dell'algoritmo di ricerca approssimata di Wu e Manber. Dire se 1011 può essere la parola D^1_i .

No. Diamo attenzione al secondo bit della parola, quello a 1.

$$D^0_i[2] = 1 \Rightarrow P[1,2] = \text{suff}_0(T[1,i]) \Rightarrow P[1,2] = \text{suff}_1(T[1,i])$$

Questo bit indica che il prefisso lungo 2 del pattern ha distanza di edit minore uguale a zero con una sottostringa del testo che finisce in posizione i .

Quindi questo è vero anche quando la distanza di edit deve essere minore uguale ad 1.

$$P[1,2] = \text{suff}_1(T[1,i]) \Rightarrow D^1_i[2] = \mathbf{1}$$

RISPOSTA: **NO**

In generale,

$$D^0_i[j] = 1 \Rightarrow D^h_i[j] = 1, h > 0$$

$$D^{h'}_i[j] = 1 \Rightarrow D^h_i[j] = 1, h > h'$$

Questa è la generalizzazione. Per qualsiasi i e j , se il bit è 1 allora se h è maggiore sarà per forza 1 anche dopo.

Esercizio 2

La parola D^0_i dell'algoritmo di ricerca approssimata di Wu e Manber è uguale a 10100. Indicare quali bit della parola D^1_i è possibile dedurre da D^0_i .

$$D^0_i = \underline{10100}$$

Innanzitutto abbiamo quanto detto nell'esercizio precedente, e abbiamo quindi questi.

In generale:

$$D^0_i[j] = 1 \Rightarrow D^1_i[j] = 1$$

Possiamo dire qualcosa sui bit vicini? (quindi quello in mezzo ai due 1, e quello subito a destra)

$$D^0_i[1] = 1 \Rightarrow D^1_i[1] = 1$$

$$D^0_i[2] = 1 \Rightarrow D^1_i[2] = 1$$

$$\Rightarrow D^1_i = \underline{1?1??}$$

$$D^0_i[1] = 1 \Rightarrow P[1,1] = \text{suff}_0(T[1,i]) \Rightarrow P[1,2] = \text{suff}_1(T[1,i])$$

$$P[1,2] = \text{suff}_1(T[1,i]) \Rightarrow D^1_i[2] = 1$$

Se analizziamo il primo 1 abbiamo che il prefisso lungo 1 del pattern ha distanza di edit al più 0 (uguale) ad una sottostringa del testo che finisce in i. Questo implica che il prefisso lungo 2 del pattern ha distanza di edit che è ≤ 1 con una sottostringa che finisce in i. Quindi il bit di posizione 2 è 1.

$$D^0_i[3] = 1 \Rightarrow P[1,3] = \text{suff}_0(T[1,i]) \Rightarrow P[1,4] = \text{suff}_1(T[1,i])$$

$$P[1,4] = \text{suff}_1(T[1,i]) \Rightarrow D^1_i[4] = 1$$

Nello stesso modo anche il bit 4 sarà 1.

In generale:

$$D^0_i[j] = 1 \Rightarrow D^1_i[j+1] = 1$$

In generale:

$$D^h_i[j] = 1 \Rightarrow D^{h+1}_i[j+1] = 1$$