

2021. november 2.

Karakter soratok hasonlósága

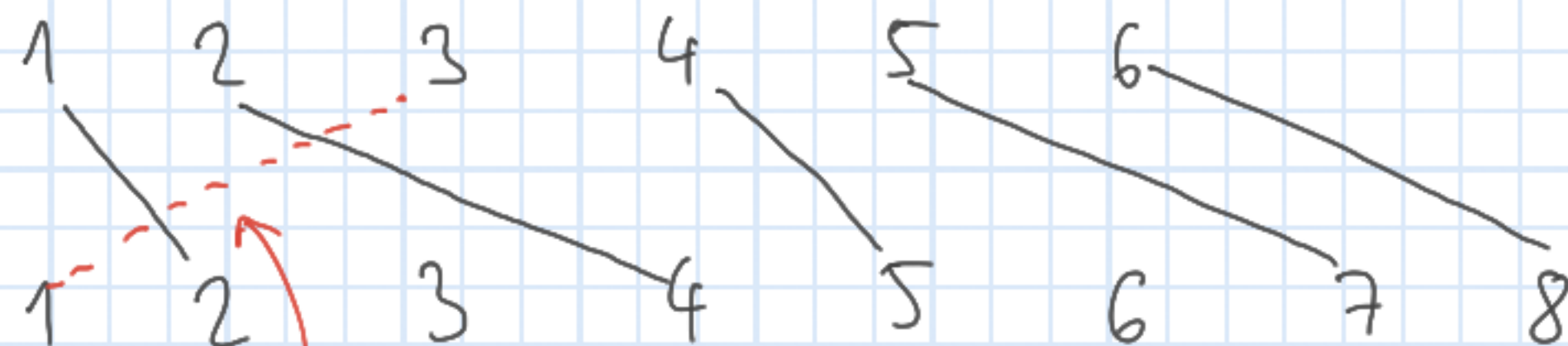
Szekvenciaillesztés (Levenshtein távolság)

Adottak az $X = (x_1, x_2, \dots, x_m)$ és $Y = (y_1, y_2, \dots, y_n)$ karakter soratok. A hasonlóságnál egy lehetséges mértéke a következő (elegendően általános lesz ez):

Teljesítsd az X és Y pozícióival $\{1, 2, \dots, m\}$ és $\{1, 2, \dots, n\}$ halmazát, és legyen M egy páros-

háa e két halmozható.

M -et rendezéscsökkentéssel rendezzük ha $(i, j) \in M$,
 $(i', j') \in M$ és $i < i'$ esetén $j < j'$



OK!

Ez nem lehet!

A rendezéscsökkentés egy vizualizációja (példa):

Legyen a két karaktersorozat:

^{1 2 3 4 5 6 7 8 9}
G A T C G G C A T

C A A T G T G A A T C
_{1 2 3 4 5 6 7 8 9 10 11}

Legyen a névenciai mentés:

$\{(1,1), (2,3), (3,4), (5,5), (6,7), (7,8), (8,9), (9,10)\}$

\Rightarrow G - A T C G - G C A T -
C A A T - G T G A A T C

Ez a megjelenítési feltehető mint az első karakter -

tersorokat áttérkeformálva a másodikba, a referencia-illesztéssel megfelelően:

cseréljük le az első G betűt C -ra, majd írjuk be egy A -t, aztán töröljük a negyedik pozíciót a C -t, és így tovább ...

Megjegyzés

Két karaktertetsorot körözt rengeteg referenciaillesztés megadható \rightarrow optimális ???

Egy M szöveciaillesztés költségét definiáljuk így:

- felhasználunk adott $g > 0$ költséget minden olyan X - és Y -beli pozícióra, amelyek nem szerepelnek M -ben
- ha $(i, j) \in M$, akkor számoljunk fel $c[x_i, y_j] \geq 0$ költséget erre a párra, ahol a $c["p", "q"]$ mennyiségek szintén adottak ($c["p", "p"] = 0$ tipikusan, de nem feltétlenül)

Adjuk össze ezeket!

Feladat

Adjunk hatékony algoritmust a minimális költségű referenciaillesztés meghatározására!

Megjegyzés

ha ez a költség kicsi, akkor a két karakter-sorozat hasonlóan futtán, ha nagy, akkor különbözően

Dinamikus programozás megoldás

① Részproblémák

$R[i, j]$ legyen az $X_i = (x_1, x_2, \dots, x_i)$ és

$Y_j = (y_1, y_2, \dots, y_j)$ kezdőszóletor (prefix) optimális rekurzióellenőrzésével megmutatása

② optimális rekonstrukcióra fel.

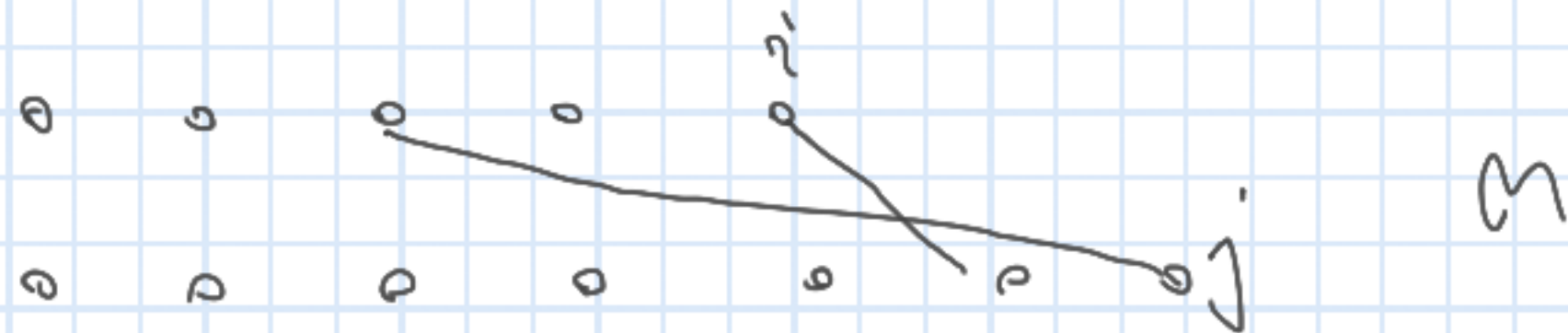
Észrevétel

Legyen M egy sorrendiajtartás X_i és Y_j között.

Ekkor három lehetőség van:

- ① $(i, j) \in M$
- ② i nem szerepel M -ben első tagként
- ③ j nem szerepel M -ben második tagként

Az is nem lehet:



Állítás

Legyen M opt. referenciaillentés X_i és Y_j között

① ha $(i, j) \in M$, akkor $M \setminus \{(i, j)\}$ opt. selv. illentés X_{i-1} és Y_{j-1} között

② ha i nem szerepel M -ben első tagként,
akkor M opt. selv. illentés X_{i-1} és Y_j között

③ ha j nem szerepel M -ben második tagként,
akkor M opt. selv. illentés X_i és Y_{j-1} között

Biz

horázos CWT & PASTE

Pl ① Ha $M \setminus \{(i, j)\}$ nem optimális, akkor
egy mála kisebb költségűtiegészítve (i, j) -vel
egy M -nél kisebb költségűt kapunk x_i -ben
és x_j -ben

③ Rekurzio

Jelölje $A[i, j]$ az $R[i, j]$ opt. megoldásának

költségét.

Nyilván $A[i, 0] = ig$ és $A[0, j] = jg$.

Ha $i, j \neq 0$, akkor a három esetben megfelelően

- $A[i, j] = A[i-1, j-1] + c[x_i, y_j]$

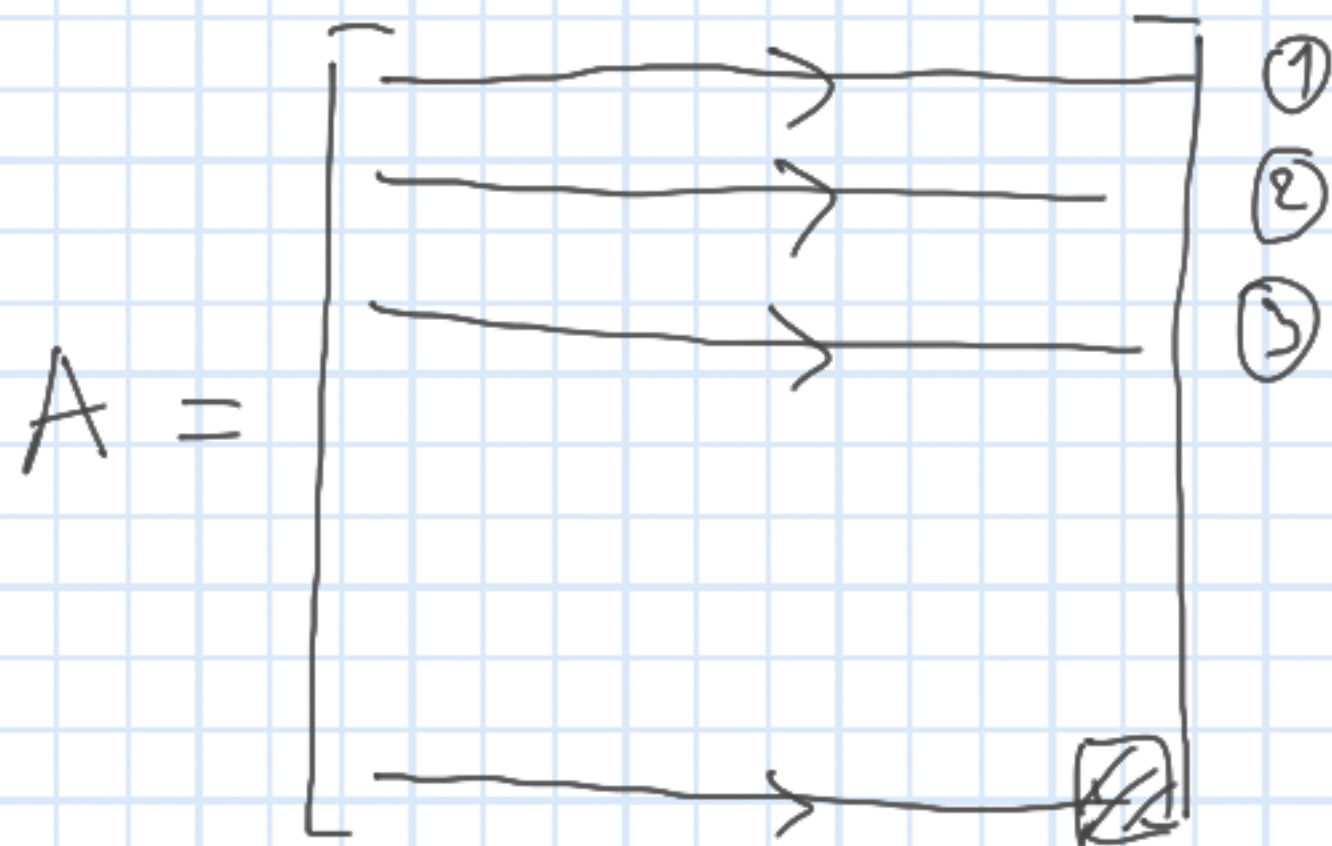
- $A[i, j] = A[i-1, j] + g$

- $A[i, j] = A[i, j-1] + g$

Nem tudjuk melyik áll fenn, így megvizsgáljuk mindehármat és a legkedvezőbbet választjuk:

$$A[i,j] = \min \begin{cases} A[i-1,j-1] + c[x_i, y_j] \\ A[i-1,j] + g \\ A[i,j-1] + g \end{cases}$$

④ optimális megoldás



$A[m,n]$ opt. megoldás

Költség : $O(mn)$ cella cellánként $O(1)$
hámoslás $\rightarrow O(mn)$

⑤ opt. relv. ill.

Minden $A[i,j]$ -nél feljegyezzük, hogy
a három esetből melyik volt a legked-
vezőbb, ebből visszafejthető az opt.
relv. ill.

