

Algorithmen und Datenstrukturen (Th. Ottmann und P. Widmayer)

Folien: Editierdistanz

Autor: Sven Schuierer

Institut für Informatik
Georges-Köhler-Allee
Albert-Ludwigs-Universität Freiburg

1 Editier-Distanz

Gegeben: zwei Zeichenketten $A = a_1 a_2 \cdots a_m$ und
 $B = b_1 b_2 \cdots b_n$

Gesucht: Minimale Kosten $D(A, B)$ für eine Folge von Editieroperationen, um A in B zu überführen

Editieroperationen:

1. **Ersetzen** eines Zeichens von A durch ein Zeichen von B
2. **Löschen** eines Zeichens von A
3. **Einfügen** eines Zeichens von B

Einheitskostenmodell:

$$c(a, b) = \begin{cases} 1 & \text{falls } a \neq b \\ 0 & \text{falls } a = b \end{cases}$$

$a = \epsilon, b = \epsilon$ möglich

Dreiecksungleichung soll für c im allgemeinen gelten:

$$c(a, c) \leq c(a, b) + c(b, c)$$

⇒ Ein Buchstabe wird höchstens einmal verändert

2 Spur

Spur als Repräsentation von Editiersequenzen

$$\begin{array}{rcccccccc} A = & & b & a & a & c & a & a & b & c \\ & & | & | & / & / & | & / & & \\ B = & a & b & a & c & b & c & a & c & \end{array}$$

oder mit

$$\begin{array}{rcccccccc} A = & - & b & a & a & c & a & - & a & b & c \\ & & | & | & & | & | & & | & & | \\ B = & a & b & a & - & c & b & c & a & - & c \end{array}$$

Kosten: 5

Aufteilung einer optimalen Spur \Rightarrow zwei optimale Teilspuren

Bemerkung: Dynamische Programmierung erforderlich!

Berechnung der Spur

Sei $A_i = a_1 \cdots a_i$ und $B_j = b_1 \cdots b_j$

$$D_{i,j} = D(A_i, B_j)$$

Drei Möglichkeiten eine Spur zu beenden:

1. a_m wird durch b_n ersetzt:

$$D_{m,n} = D_{m-1,n-1} + c(a_m, b_n)$$

2. a_m wird gelöscht: $D_{m,n} = D_{m-1,n} + 1$

3. b_n wird eingefügt: $D_{m,n} = D_{m,n-1} + 1$

3 Rekursionsgleichung

Rekursionsgleichung, falls $m, n \geq 1$:

$$D_{m,n} = \min \left\{ \begin{array}{l} D_{m-1,n-1} + c(a_m, b_n), \\ D_{m-1,n} + 1, \\ D_{m,n-1} + 1 \end{array} \right\}$$

\Rightarrow Berechnung aller $D_{i,j}$ erforderlich, $0 \leq i \leq m$,
 $0 \leq j \leq n$.

Anfangsbedingungen

Anfangsbedingungen:

$$D_{0,0} = D(\epsilon, \epsilon) = 0$$

$$D_{0,j} = D(\epsilon, B_j) = j$$

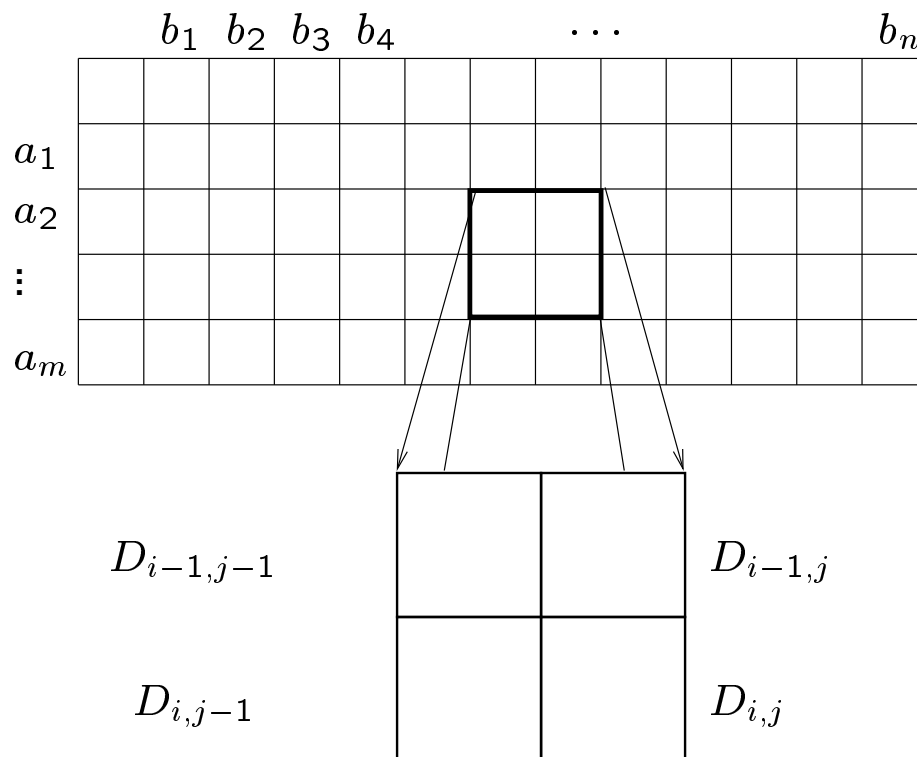
$$D_{i,0} = D(A_i, \epsilon) = i$$

Rekursionsgleichung:

$$D_{i,j} = \min \left\{ \begin{array}{l} D_{i-1,j-1} + c(a_i, b_j), \\ D_{i-1,j} + 1, \\ D_{i,j-1} + 1 \end{array} \right\}$$

4 Reihenfolge

Reihenfolge der Berechnung:



5 Algorithmus

Editierdistanz

Input: Zwei Zeichenketten $A = a_1 \cdots a_m$ und

$$B = b_1 \cdots b_n$$

Output: Die Matrix $D = (D_{ij})$

1 $D[0, 0] := 0$

2 **for** $i := 1$ **to** m **do** $D[i, 0] = i$

3 **for** $j := 1$ **to** n **do** $D[0, j] = j$

4 **for** $i := 1$ **to** m **do**

5 **for** $j := 1$ **to** n **do**

6 $D[i, j] := \min(D[i - 1, j] + 1,$
 $D[i, j - 1] + 1,$
 $D[i - 1, j - 1] + c(a_i, b_j))$

6 Beispiel

		a	b	a	c
	0	1	2	3	4
b	1				
a	2				
a	3				
c	4				

7 Editier-Operationen

Berechnung der Editieroperationen:

Editieroperationen(i, j)

Input: Berechnete Matrix D

1 if $i = 0$ and $j = 0$ then return

2 if $i \neq 0$ and $D[i, j] = D[i - 1, j] + 1$

3 then "lösche $a[i]$ "

4 Editieroperationen($i - 1, j$)

5 else if $j \neq 0$ and $D[i, j] = D[i, j - 1] + 1$

6 then "füge $b[j]$ ein"

7 Editieroperationen($i, j - 1$)

8 else

 /* $D[i, j] = D[i - 1, j - 1] + c(a[i], b[j])$ */

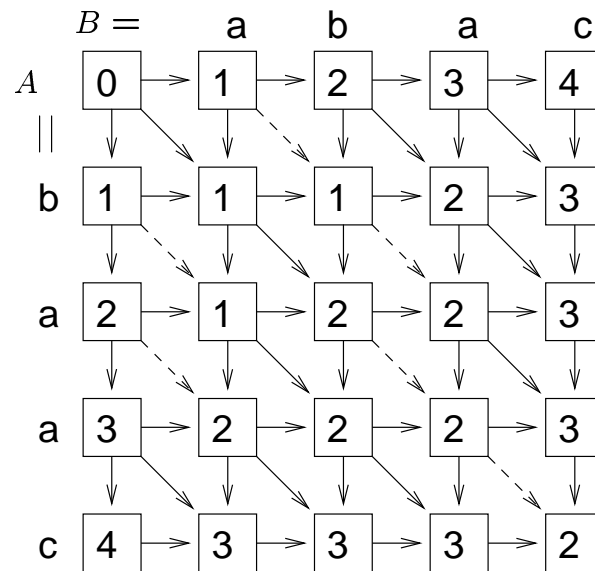
9 "ersetze $a[i]$ durch $b[j]$ "

10 Editieroperationen($i - 1, j - 1$)

Aufruf: Editieroperationen(m, n)

Beispiel

Editieroperationen:



Spurgraph: Übersicht über alle möglichen Spuren zu Transformation von A in B

Bemerkungen: Implizit vorher verwendet

Spurgraph: Knoten für jedes Paar (i, j)

gerichtete Kante von Knoten (i, j) zu (i_1, j) , $(i, j + 1)$ und $(i + 1, j + 1)$.

Gewichtung der Kanten entsprechend Editierkosten

Kosten nehmen entlang eines Weges monoton zu

Jeder Weg von der linken oberen Ecke zu rechten

unteren Ecke entspricht einer optimalen Spur

8 Approximative Zeichenkettensuche

Gegeben: zwei Zeichenketten $P = p_1p_2 \cdots p_m$ (Muster)
und $T = t_1t_2 \cdots t_n$ (Text)

Gesucht: Ein Intervall $[j', j]$, $1 \leq j' \leq j \leq n$, so daß
das Teilwort $T_{j',j} = t_{j'} \cdots t_j$ das dem Muster
 P ähnlichste Teilwort von T ist, d.h. für alle
anderen Intervalle $[k', k]$, $1 \leq k' \leq k \leq n$,
gilt:

$$D(P, T_{j',j}) \leq D(P, T_{k',k})$$

Naives Verfahren:

for all $1 \leq j' \leq j \leq n$ do

 Berechne $D(P, T_{j',j})$

 wähle Minimum

Methode

for all $1 \leq j \leq n$ do

Berechne j' , so daß $D(P, T_{j',j})$ minimal ist

Für $0 \leq i \leq m$ und $0 \leq j \leq n$ sei:

$$E_{i,j} = \min_{1 \leq j' \leq j+1} D(P_i, T_{j',j})$$

Optimale Spur:

$$\begin{array}{rcccccccc} P_i & = & b & a & a & c & a & a & b & c \\ & & | & | & / & / & | & / & & \\ T_{j',j} & = & b & a & c & b & c & a & c & \end{array}$$

Rekursionsgleichung:

$$E_{i,j} = \min \left\{ \begin{array}{l} E_{i-1,j-1} + c(p_i, t_j), \\ E_{i-1,j} + 1, \\ E_{i,j-1} + 1 \end{array} \right\}$$

Bemerkungen

Wieder drei Möglichkeiten wie die optimale Spur für $D(P_i, T_{j',j})$ endet, d.h. unabhängig von j' muß eine der drei Möglichkeiten eintreten

Aber: j' kann für $E_{i-1,j-1}$, $E_{i-1,j}$ und $E_{i,j-1}$ ganz verschieden sein.

Teilspur einer optimalen Spur ist eine optimale Teilspur, d.h. insbesondere ist j' optimal für $D(P_i, T_{j',j})$ und z.B. $E_{i,j} = E_{i,j-1} + 1$, dann ist j' auch optimal für $D(P_i, T_{j',j-1})$.

Anfangsbedingungen

$$E_{0,0} = E(\epsilon, \epsilon) = 0$$

$$E_{i,0} = E(P_i, \epsilon) = i$$

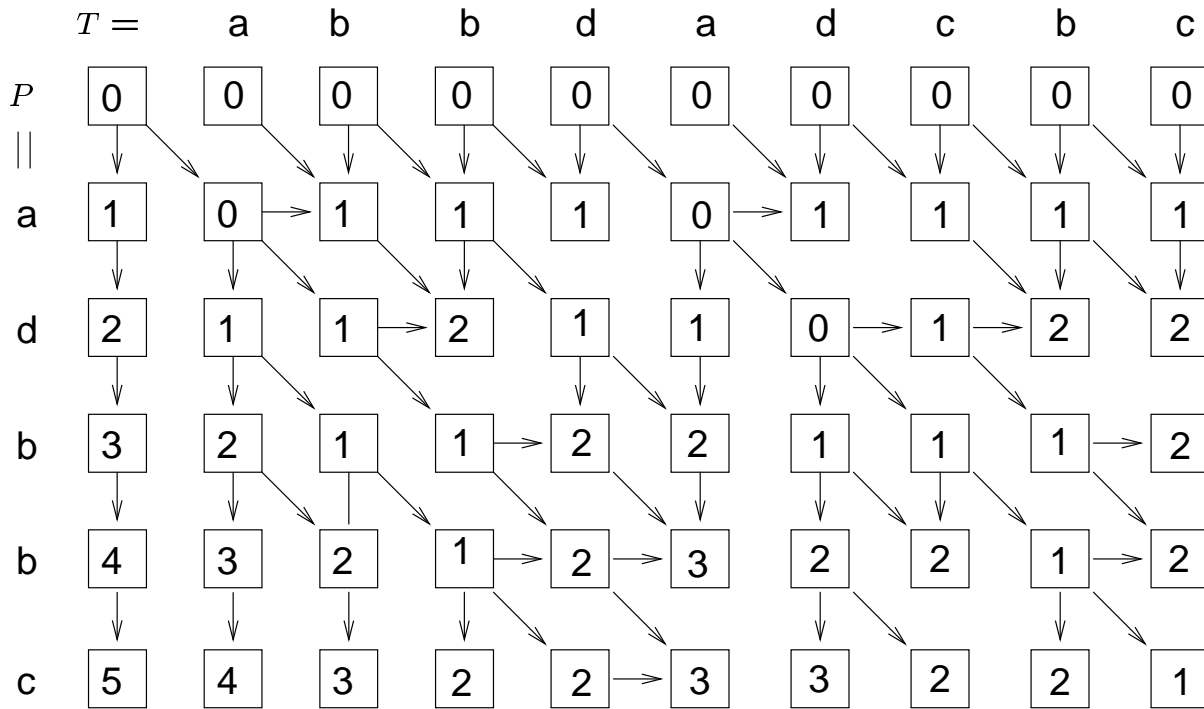
aber

$$E_{0,j} = E(\epsilon, T_j) = 0$$

Beobachtung:

Die optimale Editiersequenz von P nach $T_{j',j}$ beginnt nicht mit einer Einfügung von $t_{j'}$

Abhängigkeitsgraph



Eigentlich: Suche nach Vorkommen von P in T mit Editierdistanz höchstens k

Satz Gibt es im Abhängigkeitsgraphen einen Weg von $E_{0,j'-1}$ nach $E_{i,j}$, so ist $T_{j',j}$ ein zu P_i ähnlichstes, bei j endendes Teilstück von A mit

$$D(P_i, T_{j',j}) = E_{i,j}.$$

Ähnlichkeit zweier Zeichenketten

$S(A, B)$

Operationen:

1. Ersetzen eines Zeichens a durch ein Zeichen b :
Gewinn $s(a, b)$
2. Löschen eines Zeichens von A , Einfügen eines Zeichens von B : Verlust $-c$

Aufgabe:

Finde eine Folge von Operationen zur Umwandlung von A in B , so daß die Summe der Gewinne maximiert wird.

Motivation

Problem aus der **Biologie**: Berechnung der Ähnlichkeit von Proteinen.

Jedes **Protein** ist eine lineare Verknüpfung von 250-2000 **Aminosäuren**

Es gibt 20 Aminosäuren, aus denen alle Proteine aufgebaut sind, d.h. jedes Protein kann als **Zeichenkette über einem Alphabet** der Größe 20 aufgefaßt werden.

Manche Aminosäuren sind sich in ihren Eigenschaften (z.B. hydrophil-hydrophob) **ähnlicher** als andere

Ähnlichkeit zweier Zeichenketten nur dann sinnvoll, wenn jedes Zeichen nur in einer Operation vorkommt

Ähnlichkeit von Zeichenketten

$$S_{i,j} = S(A_i, B_j), 0 \leq i \leq m, 0 \leq j \leq n$$

Rekursionsgleichung:

$$S_{m,n} = \max \left(S_{m-1,n-1} + s(a_m, b_n), \right. \\ \left. S_{m-1,n} - c, S_{m,n-1} - c \right)$$

Anfangsbedingungen:

$$S_{0,0} = S(\epsilon, \epsilon) = 0$$

$$S_{0,j} = S(\epsilon, B_j) = -jc$$

$$S_{i,0} = S(A_i, \epsilon) = -ic$$

Ähnlichste Teilzeichenketten

Gegeben: zwei Zeichenketten $A = a_1 a_2 \cdots a_m$ und
 $B = b_1 b_2 \cdots b_n$

Gesucht: Ein zwei Intervalle $[i', i] \subseteq [1, m]$ und
 $[j', j] \subseteq [1, n]$, so daß:

$$S(A_{i',i}, B_{j',j}) \geq S(A_{k',k}, B_{l',l}),$$

für alle $[k', k] \subseteq [1, m]$ und $[l', l] \subseteq [1, n]$.

Naives Verfahren:

for all $[i', i] \subseteq [1, m]$ **and** $[j', j] \subseteq [1, n]$ **do**
 Berechne $S(A_{i',i}, B_{j',j})$

Methode

for all $1 \leq i \leq m, 1 \leq j \leq n$ do

Berechne i' und j' , so daß $S(A_{i',i}, B_{j',j})$ maximal ist

Für $0 \leq i \leq m$ und $0 \leq j \leq n$ sei:

$$H_{i,j} = \max_{\substack{1 \leq i' \leq i+1, \\ 1 \leq j' \leq j+1}} S(A_{i',i}, B_{j',j})$$

Optimale Spur:

$$\begin{array}{rcccccccc} A_{i',i} = & b & a & a & c & a & - & a & b & c \\ & | & | & & | & | & & | & & | \\ B_{j',j} = & b & a & - & c & b & c & a & - & c \end{array}$$

Rekursionsgleichung

:

$$H_{i,j} = \max \left\{ \begin{array}{l} H_{i-1,j-1} + s(a_i, b_j), \\ H_{i-1,j} - c, \\ H_{i,j-1} - c, \\ 0 \end{array} \right\}$$

Anfangsbedingungen:

$$H_{0,0} = H(\epsilon, \epsilon) = 0$$

$$H_{i,0} = H(A_i, \epsilon) = 0$$

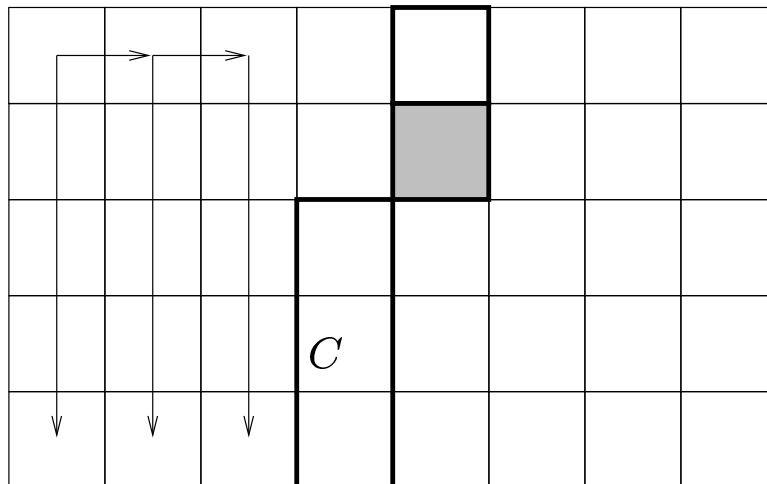
$$H_{0,j} = H(\epsilon, B_j) = 0$$

9 Editier-Distanz mit linearem Platz

Berechnung von $D_{m,n}$ mit Speicherplatz

$O(\min\{m, n\})$:

Sei $m \leq n$: Idee spaltenweise Berechnung von $D_{i,j}$



Editier-Distanz mit linearem Platz

Editierdistanz

Input: Zeichenketten A und B mit Länge m und n

```
1 for  $i := 0$  to  $m$  do  $C[j] := i$ 
2 for  $j := 1$  to  $n$  do
  /* Berechne  $D[0, j], \dots, D[m, j]$  */
  /* Inv.:  $C[i] = D[i, j - 1], i = 0, \dots, m$  */
3  $F := C[0]$  /*  $= D[0, j - 1]$  */
4  $C[0] := C[0] + 1$  /*  $= D[0, j]$  */
5 for  $i := 1$  to  $m$  do
6    $E := F$  /*  $= D[i - 1, j - 1]$  */
7    $F := C[i]$  /*  $= D[i, j - 1]$  */
8    $C[i] := \min(E + c(a_i, b_j),$ 
                 $C[i - 1] + 1, F + 1)$ 
9 return  $C[m]$ 
```