

Informatik II - SS 2016

(Algorithmen & Datenstrukturen)

Vorlesung 21 (15.7.2016)

String Matching (Textsuche)

Approximate String Matching



**UNI
FREIBURG**

Fabian Kuhn

Algorithmen und Komplexität

Gegeben:

- Zwei Zeichenketten (Strings)
- Text T (typischerweise lang)
- Muster P (engl. pattern, typischerweise kurz)

Ziel:

- Finde alle Vorkommen von P in T

Notation:

- Länge Text T : n , Länge Muster P : m

- Naiver Algorithmus
 - Gehe den Text von links nach rechts durch
 - Laufzeit $O(nm)$
- Rabin Karp Algorithmus
 - Hashe Muster und Fenster, so dass man in $O(1)$ vergleichen kann
(bei gleichem Hashwert \rightarrow ziffernweise vergleichen)
 - Beim Verschieben des Fensters wird der neue Hashwert in Zeit $O(1)$ berechnet
 - Laufzeit bei optimaler Hashfunktion $O(n+km)$

- Wir merken uns an jeder Stelle des Musters, wie weit wir das Suchfenster **sicher** bei einem “Mismatch” weiterschieben können.
- Äquivalent dazu: Stelle im Muster an der wir weiter suchen müssen
- Dies soll unabhängig vom Text sein

Vorbereitung: Array S der Länge $m + 1$

- $S[i]$: Stelle in P , an welcher man die neue Suche beginnt, falls beim Testen der Stelle i im Pattern ein Mismatch auftritt
- $S[m]$: Stelle in P , an welcher man weitersucht, nachdem P erfolgreich gefunden wurde
- $S[0] = -1, \quad S[1] = 0$

Vorberechnung von Array S :

- $P = [A, B, D, A, B, L, A, B, D, A, B, D]$
 $S = [-1, 0, 0, 0, 1, 2, 0, 1, 2, 3, 4, 5, 3]$

- An Position in $S[i]$ (für $i \in \{2, \dots, m\}$) steht

$$S[i] := \underset{k < i}{\operatorname{argmax}} \{ P[i - k \dots i - 1] = P[0 \dots k - 1] \}$$

- $S[i]$: Länge des längsten **echten** Teilstückes von $P[0 \dots i - 1]$, welches an Stelle $i - 1$ endet, und welches auch Anfangsstück von P ist
- In gewissen Fällen könnte man sicher weiter schieben (aber obiges ist ausreichend für Laufzeit $O(n)$)

Vorbereitung von Array S :

- $S[i]$: Länge des längsten *passenden* Teilstückes von $P[1 \dots i - 1]$,
- $P = [A, B, D, A, B, L, A, B, D, A, B, D]$

- $S = [-1, 0, 0, 0, 1, 2, 0, 1, 2, 3, 4, 5, 3]$

Berechnung von $S[i]$: (von links nach rechts)

- $S[0] := -1, S[1] := 0$
- Für $i > 1$ suchen wir Teilstücke, die wir verlängern können
 - Erster Kandidat: $S[i-1]$
Falls $P[S[i-1]] = P[i-1]$, dann ist $S[i] = S[i-1] + 1$
 - Zweiter Kandidat: $S[S[i-1]]$
Falls $P[S[S[i-1]]] = P[i-1]$, dann ist $S[i] = S[S[i-1]] + 1$
 - Usw...

Berechnung von $S[i]$: Beispiel

$h := S[i - 1]$

while $h \geq 0$ do

 if $P[i - 1] = P[h]$ then

$S[i] := h + 1; h := -1$

 else

$h := S[h]$

if $h = -1$ then $S[i] = 0$

$h := S[i - 1]$

while $h \geq 0$ do

 if $P[i - 1] = P[h]$ then

$S[i] := h + 1; h := -1$

 else

$h := S[h]$

if $h = -1$ then $S[i] = 0$

Beobachtung:

$$S[i] \leq S[i - 1] + 1$$

Falls $S[i] = S[i - 1] + 1$: 1 Schleifendurchlauf

Falls $S[i] < S[i - 1]$:

- Wert von h nimmt in jedem Schleifendurchlauf ab
- Am Schluss ist $S[i] = h + 1$
- #Schleifendurchläufe $\leq \Delta h + 1 \leq S[i - 1] - (S[i] - 1) + 1$
 $\leq S[i - 1] - (S[i]) + 2$

Falls $S[i] = S[i - 1] + 1$:

- Anzahl Schleifendurchläufe = $1 = S[i - 1] - S[i] + 2$

Falls $S[i] < S[i - 1]$:

- Anzahl Schleifendurchläufe $\leq \Delta h + 1 = S[i - 1] - S[i] + 2$

Gesamtlaufzeit:

Knuth-Morris-Pratt Algorithmus:

- Berechnet zuerst in Zeit $O(m)$ das Array S der Länge m
 - hängt nur vom Pattern P ab und nicht vom Text T
 - beschreibt an jeder Position im Pattern, wo (im Pattern) man bei einem Mismatch sicher weitersuchen kann
- Mit Hilfe von S werden dann alle Vorkommen von P in T in Zeit $O(n)$ gefunden
 - In jedem Schritt kann man entweder die aktuelle Suchposition in T oder die Position des Suchfensters in T um mindestens 1 nach rechts verschieben

Gesamtlaufzeit: $O(m + n) = O(n)$

Gegeben: Zwei Strings

$$A = a_1 a_2 \dots a_m \text{ and } B = b_1 b_2 \dots b_n$$

Ziel: Bestimme die minimale Anzahl $ED(A,B)$ der Editieroperationen, die benötigt werden um A in B zu transformieren

Editieroperationen:

- a) **Replace** a character from string A by a character from B
- b) **Delete** a character from string A
- c) **Insert** a character from string B into A

m a - t h e m - - a t i c i a n
m u l t i p l i c a t i o - - n

- Verschiedene Kosten für verschiedene Buchstaben:
 - Kosten für **replace(a,b) = $c(a, b) \geq 0$**
 - Kosten für **delete(a) = $c(a, \varepsilon)$**
 - Kosten für **insert(b) = $c(\varepsilon, b)$**

ε =leere Wort

- **Wichtig dabei: Dreiecksungleichung:**

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

(damit jeder Buchstabe maximal einmal geändert wird)

- **Unit cost model:** $c(a, b) = \begin{cases} 1, & \text{if } a \neq b \\ 0, & \text{if } a = b \end{cases}$

- $D_{k,l} = ED(a_1 a_2 \dots a_k, b_1 b_2 \dots b_l)$

- **Rekursionsgleichung** (für $k, \ell \geq 1$)

$$D_{k,\ell} = \min \left\{ \begin{array}{l} D_{k-1,\ell-1} + c(a_k, b_\ell) \\ D_{k-1,\ell} + c(a_k, \varepsilon) \\ D_{k,\ell-1} + c(\varepsilon, b_\ell) \end{array} \right\}$$

- **Basisfälle**

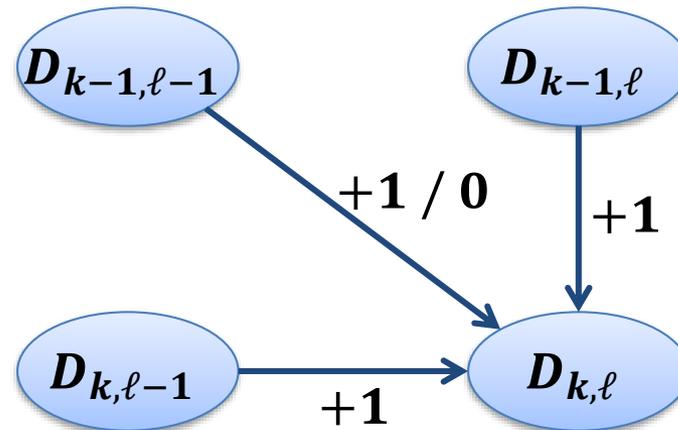
$$D_{0,0} = D(\varepsilon, \varepsilon) = 0$$

$$D_{0,j} = D(\varepsilon, B_j) = D_{0,j-1} + c(\varepsilon, b_j)$$

$$D_{i,0} = D(A_i, \varepsilon) = D_{i-1,0} + c(a_i, \varepsilon)$$

Editierdistanz – Dynamisches Programm

- Berechne $D_{i,j}$ für alle $0 \leq i \leq k, 0 \leq j \leq \ell$:



		<i>a</i>	<i>b</i>	<i>c</i>	<i>c</i>	<i>a</i>
<i>b</i>						
<i>a</i>						
<i>b</i>						
<i>d</i>						
<i>a</i>						

		<i>a</i>	<i>b</i>	<i>c</i>	<i>c</i>	<i>a</i>
	0	1	2	3	4	5
<i>b</i>	1	1	1	2	3	4
<i>a</i>	2	1	2	2	3	3
<i>b</i>	3	2	1	2	3	4
<i>d</i>	4	3	2	2	3	4
<i>a</i>	5	4	3	3	3	3

- Editierdistanz zwischen zwei Strings der Länge m und n kann in Zeit $O(mn)$ berechnet werden.
- Durch Speichern der *Vorgänger* kann man auch die zugehörigen Editieroperationen erhalten
- Erweiterung auf komplexere Kostenmodelle ist ohne Probleme möglich

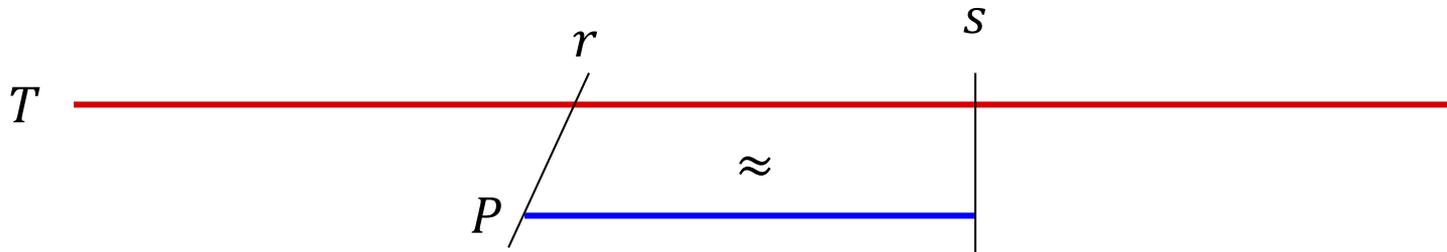
Gegeben: Strings

$T = t_1 t_2 \dots t_n$ (Text) und

$P = p_1 p_2 \dots p_m$ (Muster).

Ziel: Finde ein Intervall $[r, s]$, $1 \leq r \leq s \leq n$ so dass der Teilstring $T_{r,s} := t_r \dots t_s$ die größte Ähnlichkeit zum Muster P hat:

$$\arg \min_{1 \leq r \leq s \leq n} D(T_{r,s}, P)$$



Naive Lösung:

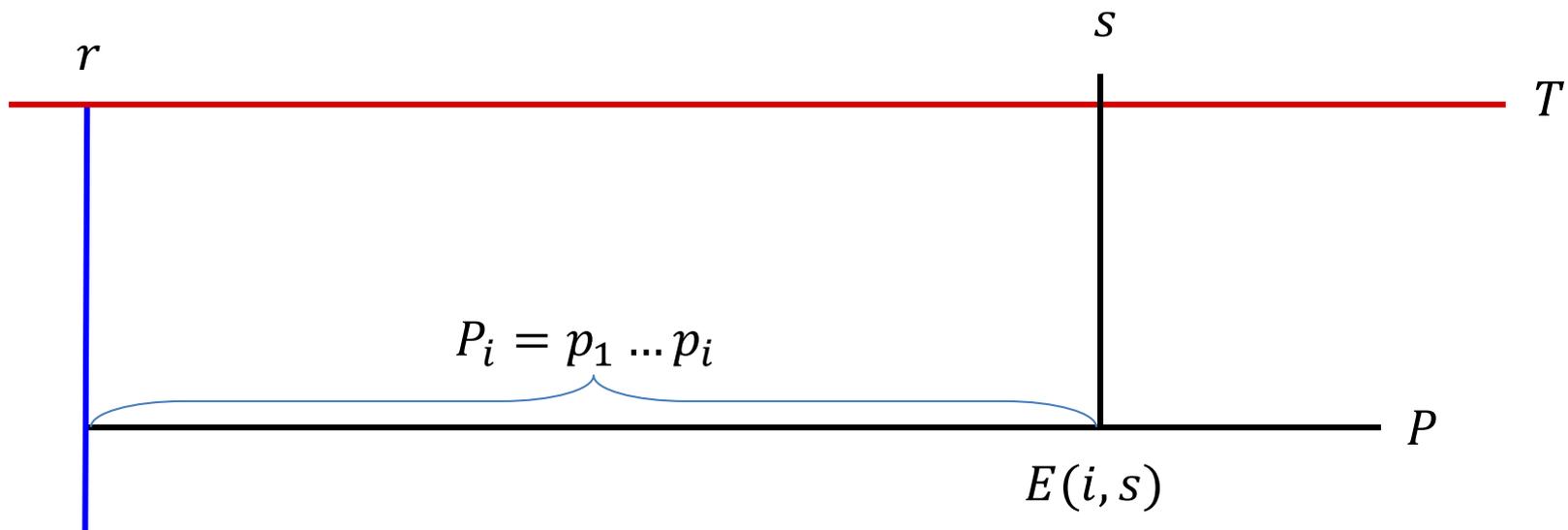
for all $1 \leq r \leq s \leq n$ do

 compute $ED(T_{r,s}, P)$

choose the minimum

Ein verwandtes Problem:

- Berechne für jede Position s im Text und i im Muster die kleinste Editierdistanz $E(i, s)$ zwischen $P_i = p_1 \dots p_i$ und einem beliebigen Teilstring $T_{r,s}$ von T , der in Position s endet.



Es gibt drei Möglichkeiten ein *Alignment* von T_s and P_i zu beenden:

1. t_s is replaced by p_i :

$$E_{s,i} = E_{s-1,i-1} + c(t_s, p_i)$$

2. t_s is deleted:

$$E_{s,i} = E_{s-1,i} + c(t_s, \varepsilon)$$

3. p_i is inserted:

$$E_{s,i} = E_{s,i-1} + c(\varepsilon, p_i)$$

$E_{i,j}$ = Editierdistanz eines optimalen Matchings von P_i
welches an Position j im Text endet.

Basisfälle:

$$E_{0,0} = 0$$

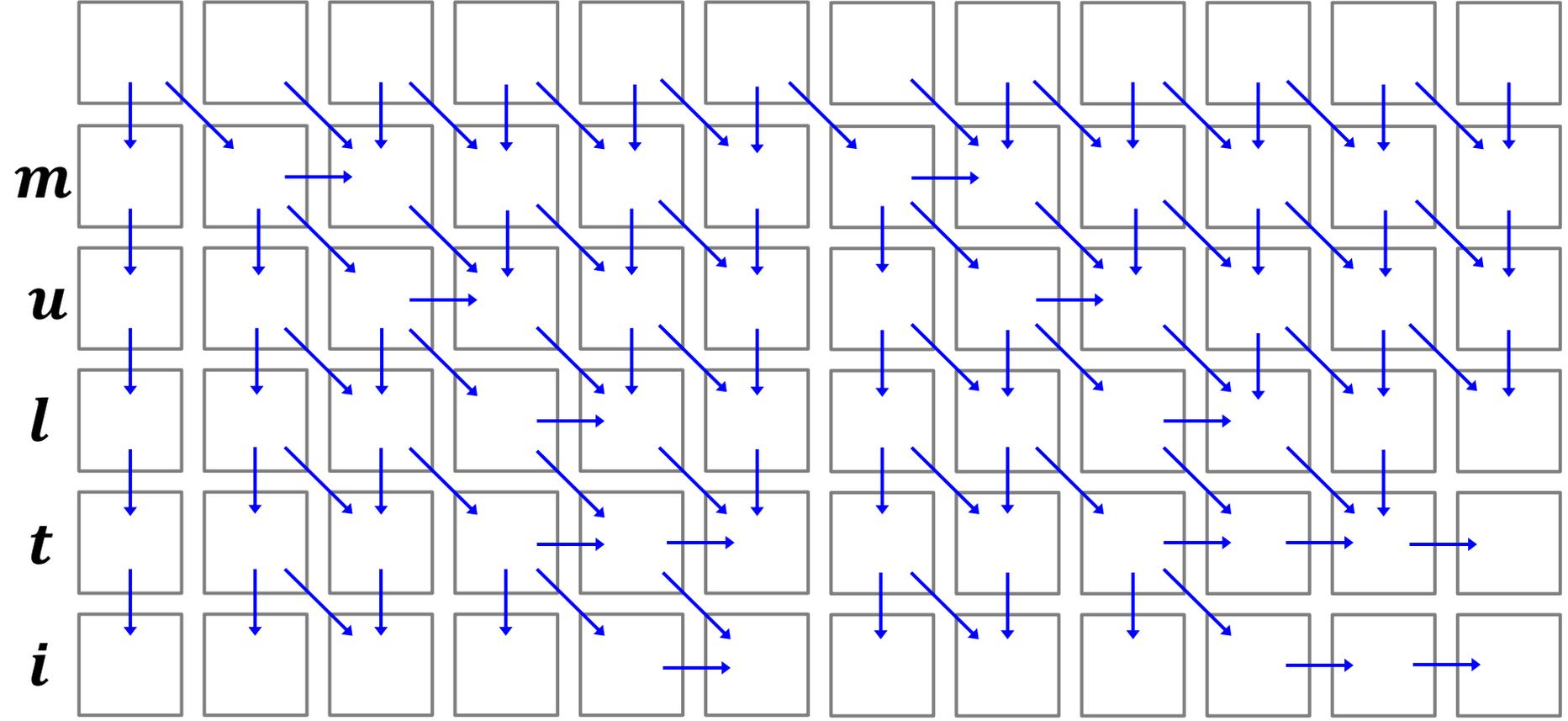
$$E_{0,j} = j$$

$$E_{i,0} = 0$$

Rekursionsgleichung (unit cost model):

$$E_{s,i} = \min \left\{ \begin{array}{l} E_{s-1,i-1} + 1/0 \\ E_{s-1,i} + 1 \\ E_{s,i-1} + 1 \end{array} \right\}$$

m a t h e m a t i c s



- Algorithmus um $E(m, n)$ zu berechnen ist bis auf die Initialisierung von $E(i, 0)$ identisch zum Editierdistanzalgorithmus
- Bei der Editierdistanz: Nur eintrag $E(m, n)$ zählt
- Approximate String Matching: Minimum der letzten Zeile zählt.
- Approximate String Matching ist in $O(nm)$ Zeit lösbar