

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 $O(n+km)$

Knuth-Morris-Pratt Alg.: Initialisierung

- 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 (*wichtig*)

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 :

- $P = [\underline{A, B, D}, A, B, L, \underline{A, B, D}, \bar{A}, B, D]$
 $S = [-1, 0, 0, 0, 1, 2, 0, 1, 2, \boxed{3}, 4, 5, 3]$

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

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

- $S[i]$: Länge des längsten **echten** Teilstückes von $P[0 \dots i - 1]$, welches an Stelle $i - 1$ endet, und welches auch Anfangsstück von P ist
Stück, welches vor der aktuellen Position endet
- 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 $P[1 \dots i - 1]$,

- $P = [\text{A, B, D, A, B, L, A, B, D, A, B, D}]$

The diagram shows the string $P = [\text{A, B, D, A, B, L, A, B, D, A, B, D}]$. Red underlines and arrows indicate matching substrings: A, B, D (indices 1-3), A, B (indices 4-5), A, B, D (indices 7-9), and A, B (indices 10-11). Blue annotations include question marks, 'Nein', 'Ja', and the number 3.

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

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: $S[i-1] = 3$
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

$h := S[i - 1]$

// erste Kandidat zur Verlängerung

while $h \geq 0$ do

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

} erfolgreiche Verlängerung

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

else

// nächste Kandidat zur Verlängerung

$h := S[h]$

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

$h = \#$ Zeichen des Kandidaten, den wir verlängern wollen.

Berechnung von $S[i]$: Laufzeit

```
h := S[i - 1]
while h ≥ 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 $\leq \Delta h + 1 \leq S[i - 1] - (S[i] - 1) + 1$
 $\leq S[i - 1] - (S[i]) + 2$

Anfangs h
h am Ende

Berechnung von $S[i]$: Laufzeit

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

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

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



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

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

Gesamtlaufzeit:

$$\sum_{i=2}^m (S[i-1] - S[i] + 2) = 2 \cdot (m-1) + (S[1] - S[2] + S[2] - S[3] + \dots + S[m-1] - S[m])$$

$$= 2 \cdot (m-1) + \underbrace{S[1]}_{=0} - \underbrace{S[m]}_{\geq 0} \leq 2(m-1)$$

$$\leadsto O(m)$$

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

Recap - Editierdistanz

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

insert *mit Änderung* *Replace* *delete()*

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) = $c(a, b) \geq 0$** (*in $c(a, b)$ darf anders sein als $c(a', b')$*)
 - 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 $k, l \geq 1$)

$$D_{k,l} = \min \left\{ \begin{array}{l} D_{k-1,l-1} + c(a_k, b_l) \\ D_{k-1,l} + c(a_k, \varepsilon) \\ D_{k,l-1} + c(\varepsilon, b_l) \end{array} \right\} = \min \left\{ \begin{array}{l} D_{k-1,l-1} + 1/0 \\ D_{k-1,l} + 1 \\ D_{k,l-1} + 1 \end{array} \right\}$$

*nach der letzten Operation
unterscheiden*

= 0 unit cost modell

- **Basisfälle**

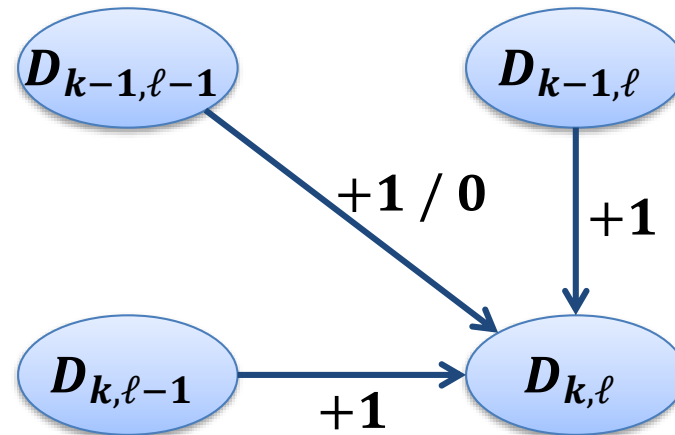
$$D_{0,0} = D(\varepsilon, \varepsilon) = 0$$

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

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

Editierdistanz – Dynamisches Programm

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



	ϵ	<u>a</u>	b	c	c	a
ϵ	0	1	2	3	4	5
<u>b</u>	1	<u>1</u>				
a	2					
b	(3)					
d	4					
a	5					

		<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

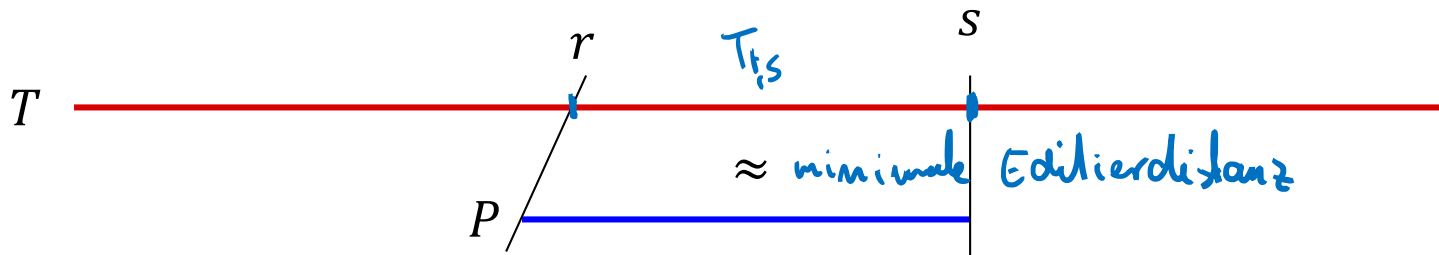
Gegeben: Strings

$T = t_1 t_2 \dots t_n$ (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 $ED(T_{r,s}, P)$ $\mathcal{O}((r-s+1) \cdot m) = \mathcal{O}(n \cdot m)$

choose the minimum

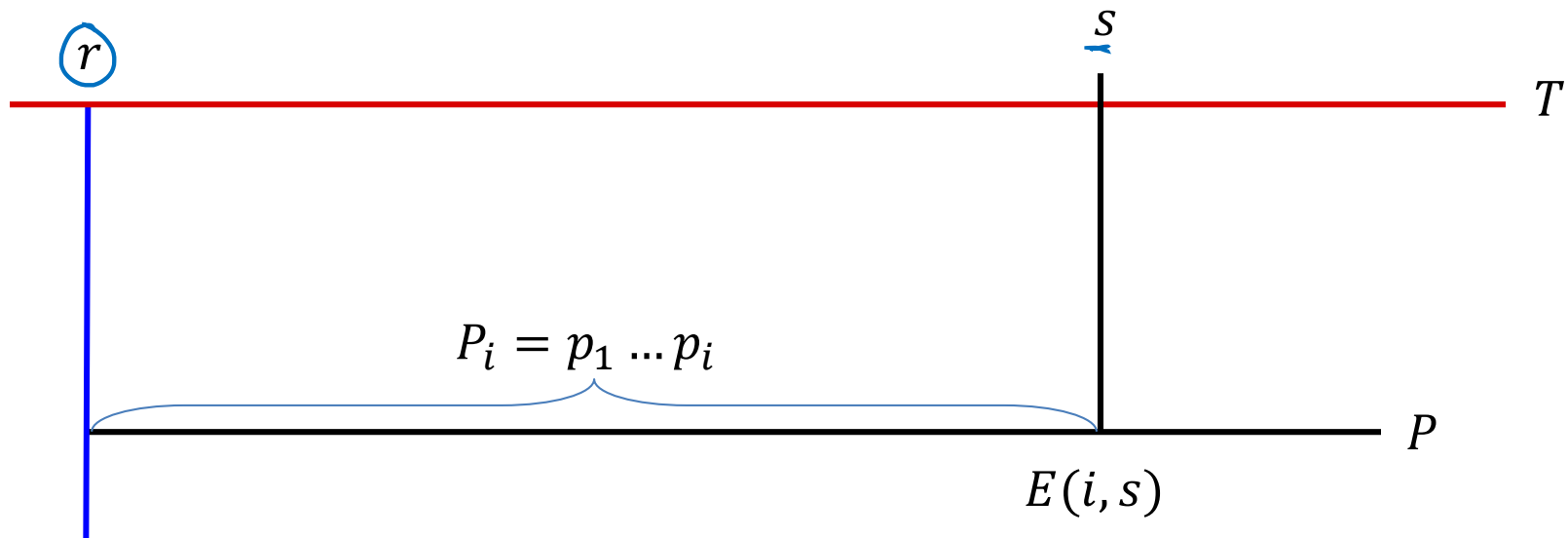
Gesamtlaufzeit ist hier

$$\mathcal{O}(n^2 \cdot n \cdot m) = \mathcal{O}(n^3 \cdot m)$$

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 / oder das leere Wort ist.
ist unbekannt

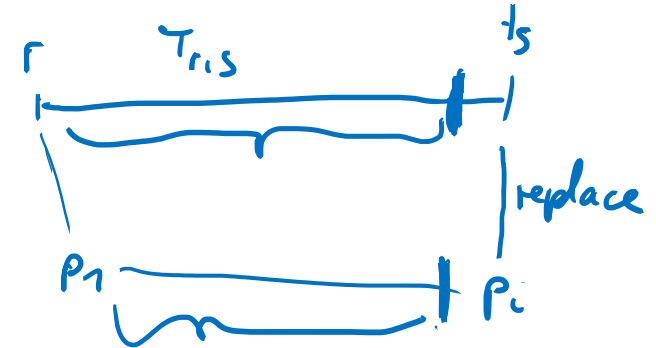


Approximate String Matching

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

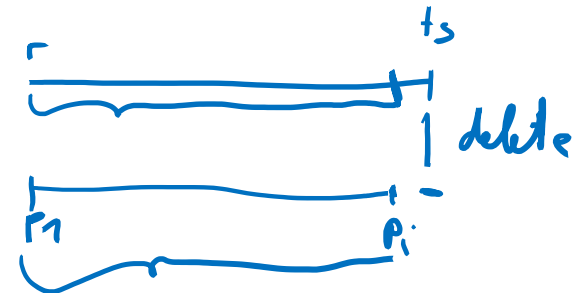
1. t_s is replaced by p_i :

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



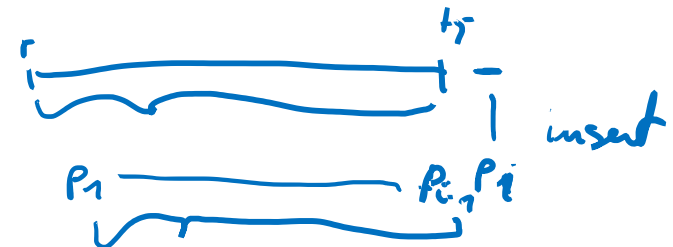
2. t_s is deleted:

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



3. p_i is inserted:

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



Approximate String Matching

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

Basisfälle:

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

$$E_{s,0} = j \quad \cancel{E_{0,j} = 1} \quad \curvearrowright$$

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

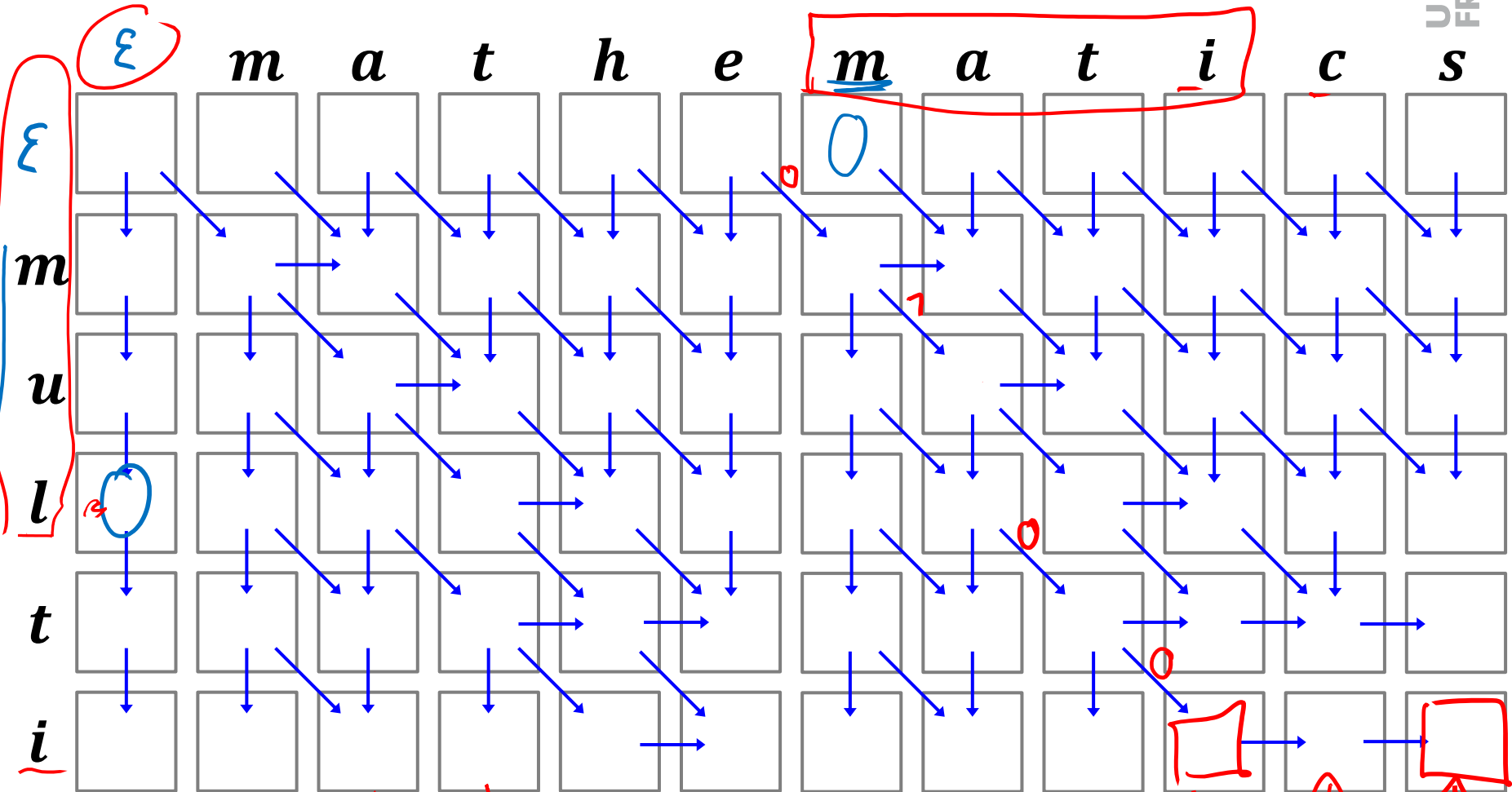
leere Wort approximativ auf das leere Wort match
Wort der Länge j approximativ auf das leere Wort match
leere Wort auf ein Wort der Länge i matchen ~~to~~ wollen

Rekursionsgleichung (unit cost model):

$$E_{s,i} = \min \left\{ \begin{array}{l} E_{s-1,i-1} + 1/0 \\ E_{s-1,i} + 1 \\ E_{s,i-1} + 1 \end{array} \right\} \quad (\text{identisch zur Editierdistanz})$$

$$E_{3,0} = 3$$

m a - E i
m u l t i



Bei der Editierdistanz: nur der Eintrag $ED(m, n)$ wichtig

App. String Matching: Suche Minimum der letzten Zeile suchen

- 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