

Informatik II: Algorithmen und Datenstrukturen SS 2013

Vorlesung 11a, Dienstag, 2. Juli 2013
(Editierdistanz, rekursive Berechnung)

Prof. Dr. Hannah Bast
Lehrstuhl für Algorithmen und Datenstrukturen
Institut für Informatik
Universität Freiburg

Blick über die Vorlesung heute

■ Organisatorisches

- Ihre Erfahrungen mit dem Ü10 (BFS, Dijkstra)

■ Editierdistanz

- Maß für die Ähnlichkeit zwischen zwei Wörtern
- Algorithmus 1: rekursiv
- Algorithmus 2: mittels "dynamischer Programmierung"
- Ü11, Aufgabe 1: Algorithmus 2 implementieren

Erfahrungen mit dem Ü10 (BFS, Dijkstra)

■ Zusammenfassung / Auszüge

Stand 2. Juli 16:00

- Test mit den "Zwischenrunden" wegen mehrfach eingefügter Knoten (ohne `changeKey`) war verwirrend ... **einverstanden**
- Dijkstra war etwas dürftig erklärt ... **hmm ?**
- Leider keine gescheite Definitionen + zu wenig Zeit
- Könnte auch gerne noch mal etwas schwerer werden
- Keine Motivation mehr, da **50%** der Punkte schon erreicht
- Danke für den Hinweis mit dem Gottesbeweis
- Frühlingspunkt: steht unter Schock / nicht an der TF
- Gource-Video für unser SVN ?

Musterlösung Ü10, Aufgabe 2

gemessen
in Kanten-
anzahl.

■ BFS berechnet kürzeste Wege mit Kantenkosten 1

- Genauer: Level $i = \{ \text{Knoten } v \text{ mit } \text{dist}(s, v) = i \}$
- Beweis durch vollständige Induktion über i
- Induktionsanfang $i = 0$: Level $0 = \{ s \}$, und das ist auch der einzige Knoten v mit $\text{dist}(s, v) = 0$
- Induktionsschritt $i \rightarrow i + 1$: z.z. Level $i+1 = \{ v : \text{dist}(s, v) = i+1 \}$
IV: $\forall j \leq i : \text{Level } j = \{ v : \text{dist}(s, v) = j \}$

gleichheit
von Mengen
also \subseteq und \supseteq

" \subseteq " : v in Level $i+1$ von BFS $\Rightarrow v$ über eine Kante mit Knoten w aus Level i verbunden.

IV $\Rightarrow \text{dist}(s, w) = i$
 $\Rightarrow \text{dist}(s, v) \leq i+1$... = folgt noch nicht!

Annahme $\text{dist}(s, v) = k \leq i$
IV $\Rightarrow v \in \text{Level } k \leq i \downarrow v$ in Level $i+1$

$\Rightarrow \text{dist}(s, v) = i+1$ \square " \supseteq " : ...

- Viele Anwendungen, wo man ähnliche Strings sucht
 - Dubletten in Adressdatenbanken
 - Hein Blöd, 27568 Bremerhaven
 - Hein Bloed, 27568 Bremerhafen
 - Hein Doof, 27478 Cuxhaven
 - Produktsuche
 - Memory Stik
 - Websuche
 - eyjaföllajaküll
 - uniwersität verien 2013
 - Gemeinsamkeit: man braucht ein Maß für die **Ehnlichkeit** zwischen zwei Strings

Editierdistanz 1/7

das zeigt erst mal
nur $ED(x, y) \leq 4$
 $x = \text{DOOF}$, $y = \text{BLOED}$

■ Definition Editierdistanz, auch Levenshtein-Distanz

- Gegeben zwei Zeichenketten (strings) x und y
- $ED(x, y)$ = Editierdistanz (edit distance) von x und y = die minimale Anzahl Operationen um x in y zu transformieren:

- Einfügen eines Buchstabens (**insert**)
- Ersetzen eines Buchstabens durch einen anderen (**replace**)
- Löschen eines Buchstabens (**delete**)
- Die **Position** einer Operation ist ... siehe Beispiel:

1 2 3 4
D O O F
B O O F
B L O F
B L O E F
B L O E D
1 2 3 4 5

REPLACE(1, B)
REPLACE(2, L)
INSERT(4, E)
REPLACE(5, D)

↑
man hat eine
Folge!

1 2 3 4 5
B L O E D
B L O E F
B L O F
B O O F
D O O F
1 2 3 4

REPLACE(5, F)
DELETE(4)
REPLACE(2, O)
REPLACE(1, D)

↑
nicht man hat

Editierdistanz 2/7

■ Etwas Notation

- Mit ε bezeichnen wir das leere Wort
- Mit $|x|$ bezeichnen wir die Länge von x (= Anzahl Zeichen)
- Mit $x[i..j]$ bezeichnen wir die Teilfolge der Zeichen i bis j der Zeichenkette x , wobei $1 \leq i \leq j \leq |x|$

■ Ein paar einfache Eigenschaften

- $ED(x, y) = ED(y, x)$ *man kann die Folge von Operationen "umkehren" ... siehe Folie 6*
- $ED(x, \varepsilon) = |x|$
- $ED(x, y) \geq \text{abs}(|x| - |y|)$ *aber nicht unbedingt =*
 $\text{abs}(x) = x > 0 ? x : -x$
- $ED(x, y) \leq ED(x[1..n-1], y[1..m-1]) + 1$ $n = |x|, m = |y|$



Editierdistanz 3/7

■ Lösungsideen anhand von Beispielen

- Für VERIEN → FERIEN ? $ED = 1$
- Für MEXIKO → AMERIKA ? $ED = 3$
- Für AAEBEAABEAREEEAEB → RBEAAEEBAAAEBBAEAE ?
- **Beobachtung:** möglichst große gemeinsame Teilstrings zu finden klappt manchmal, aber nicht immer

■ Rekursiver Ansatz

- In zwei Teile / Hälften teilen? Keine gute Idee, z.B.
 $ED(\text{GRAU}, \text{RAUM}) = 2$ aber $ED(\text{GR}, \text{RA}) + ED(\text{AU}, \text{UM}) = 4$
- Auf ein "kleineres" Problem zurückführen?

Das probieren wir jetzt !

Editierdistanz 4/7

1 2 3 4 5 6 7
S A U D O O F
A U D O O F
U D O O F
D O O F
1 2 3 4

DELETE(1)
DELETE(1)
DELETE(1)

■ Terminologie

- Seien x und y unsere beiden Zeichenketten
- Seien $\sigma_1, \dots, \sigma_k$ eine Folge von $k = ED(x, y)$ Operationen für $x \rightarrow y$, das heißt um x in y zu überführen
(Wir nehmen im Folgenden nicht an, dass wir die Folge schon kennen, sondern nur, dass es so eine gibt)
- Wir betrachten im Folgenden nur **monotone** Op.-Folgen, d.h. die Position von σ_{i+1} ist \geq die Position von σ_i , wobei = nur dann erlaubt ist, wenn beides delete Operationen sind
- **Lemma:** Für beliebige x und y mit $k = ED(x, y)$ gibt es eine **monotone** Folge von k Operationen für $x \rightarrow y$
- **Beweisintuition:** die Reihenfolge der Operationen ist im Grunde egal, also kann man sie auch monoton anordnen

siehe Folge 6

Editierdistanz 5/7

FALL 1A: $\overset{x}{D} \overset{2-1}{O} \overset{2}{O} \overset{\sigma_k}{F} \rightarrow \overset{2}{B} \overset{2}{L} \overset{2}{O} \overset{2}{E} \rightarrow \overset{2}{B} \overset{2}{L} \overset{2}{O} \overset{2}{E} \overset{2}{D}$
FALL 1B: $\overset{x}{B} \overset{2-1}{L} \overset{2}{O} \overset{2}{E} \overset{\sigma_k}{D} \rightarrow \overset{2}{D} \overset{2}{O} \overset{2}{O} \overset{2}{F} \rightarrow \overset{2}{D} \overset{2}{O} \overset{2}{O} \overset{2}{F}$
FALL 1C: $\overset{x}{D} \overset{2-1}{O} \overset{2}{O} \overset{2}{F} \overset{\sigma_k}{F} \rightarrow \overset{2}{B} \overset{2}{L} \overset{2}{O} \overset{2}{E} \overset{2}{F} \rightarrow \overset{2}{B} \overset{2}{L} \overset{2}{O} \overset{2}{E} \overset{2}{D}$
FALL 2: $\overset{x}{D} \overset{2-1}{O} \overset{2}{O} \overset{2}{F} \overset{\sigma_k}{I} \rightarrow \overset{2}{B} \overset{2}{L} \overset{2}{O} \overset{2}{E} \overset{2}{F} \overset{\sigma_k}{I} \rightarrow \overset{2}{B} \overset{2}{L} \overset{2}{O} \overset{2}{E} \overset{2}{D}$

■ Fallunterscheidung

- Wir betrachten die letzte Operation σ_k
 - $\sigma_1, \dots, \sigma_{k-1} : X \rightarrow Z$ und $\sigma_k : Z \rightarrow Y$
DOOF BLOEF BLOEF BLOED
 - Seien $n = |x|$ und $m = |y|$ und $m' = |z|$
4 5 5
 - Man beachte, dass $m' \in \{m - 1, m, m + 1\}$
- Fall 1: σ_k macht etwas "ganz am Ende" von z , d.h. eins von:
 - Fall 1a: $\sigma_k = \text{insert}(m' + 1, y[m])$ [dann ist $m' = m - 1$]
 - Fall 1b: $\sigma_k = \text{delete}(m')$ [dann ist $m' = m + 1$]
 - Fall 1c: $\sigma_k = \text{replace}(m', y[m])$ [dann ist $m' = m$]
- Fall 2: σ_k macht nichts "ganz am Ende" von z
 - dann $z[m'] = y[m]$ und $x[n] = z[m']$ und damit
 $\sigma_1, \dots, \sigma_k : x[1..n-1] \rightarrow y[1..m-1]$ und $x[n] = y[m]$

- Wir haben also einen dieser vier Fälle

- Fall 1a (**insert** am Ende): $\sigma_1, \dots, \sigma_{k-1} : x[1..n] \rightarrow y[1..m-1]$
- Fall 1b (**delete** am Ende): $\sigma_1, \dots, \sigma_{k-1} : x[1..n-1] \rightarrow y[1..m]$
- Fall 1c (**replace** am Ende): $\sigma_1, \dots, \sigma_{k-1} : x[1..n-1] \rightarrow y[1..m-1]$
- Fall 2 (**nichts** am Ende): $\sigma_1, \dots, \sigma_k : x[1..n-1] \rightarrow y[1..m-1]$

■ Daraus folgt die folgende rekursive Formel

– Für $|x| > 0$ und $|y| > 0$ ist:

$ED(x, y)$ = das Minimum von

- $ED(x[1..n], y[1..m-1]) + 1$, *insert* und *Fall 1A*
- $ED(x[1..n-1], y[1..m]) + 1$, *delete* und *Fall 1B*
- $ED(x[1..n-1], y[1..m-1]) + 1$ *replace* falls $x[n] \neq y[m]$ *Fall 1C*
- $ED(x[1..n-1], y[1..m-1])$ falls $x[n] = y[m]$ *Fall 2*

– Für $|x| = 0$ ist $ED(\overset{z}{x}, y) = |y|$

– Für $|y| = 0$ ist $ED(x, \overset{z}{y}) = |x|$

- Editierdistanz

- In Wikipedia (Definition + Algorithmen)

- http://en.wikipedia.org/wiki/Levenshtein_distance

- <http://de.wikipedia.org/wiki/Levenshtein-Distanz>