# Informatik II: Algorithmen und Datenstrukturen SS 2017

Vorlesung 11a, Dienstag, 11. Juli 2017 (Editierdistanz, Teil 1)

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

## Blick über die Vorlesung heute

UNI FREIBURG

#### Organisatorisches

Erfahrungen ÜB10

Routenplanung in BW / SL

#### ■ Inhalt

EditierdistanzMotivation + Notation

Rekursive Berechnung monotone Folgen

Rekursive Implementierung
 Code + Laufzeit

ÜB11: diesmal wieder ein reines Theorieblatt







## Erfahrungen ÜB10 1/2

#### Zusammenfassung

- Schöne Aufgabe mit Praxisbezug, hat vielen Spaß gemacht
- Einige Kämpfe mit dem Design / Code / Debuggen
   Extra weniger Vorgaben diesmal, ist ja schon das ÜB10
   Es hat sich aber gezeigt, dass viele noch Probleme mit der Umsetzung eines abstrakten Algorithmus in Code haben
- Einigen ist offenbar nicht klar, wie Sie richtig testen sollen
   Jede Funktion einzeln, siehe 10 Gebote auf dem ÜB10
   Grundsätzlich auf kleinen Beispielen, das geht immer
- Einige haben etwas Zeit gebraucht, um zu verstehen, was der "längste kürzeste Weg" ist

## Erfahrungen ÜB10 2/2

#### Ergebnis

- Erklärung längster kürzester Pfad ab einem Startknoten s
   Der Zielknoten v, für den dist(s, v) maximal ist
   Wobei dist(s, v) der kürzeste Weg von s nach v ist
- Längster kürzester Weg ab der Freiburger TF in BaWü:
   <a href="http://share.mapbbcode.org/cusya">http://share.mapbbcode.org/cusya</a> "Auf zur Kunigundenkapelle"
- Längster kürzester Weg ab der Informatik im Saarland:
   <a href="http://share.mapbbcode.org/psdbh">http://share.mapbbcode.org/psdbh</a> eine Kreuzung, laangweilig

#### Editierdistanz 1/9



#### Motivation

 Es gibt viele Anwendungen, wo man ein Maß für die Ähnlichkeit zwischen zwei Zeichenketten braucht

Drei Beispiele dazu auf der nächsten Folie

Die Editierdistanz ist ein solches Ähnlichkeitsmaß

Das in der Praxis am häufigsten verwendete Maß

Definition auf Folie 7

Benannt nach dem **Edi-Tier** 



## \_

## Editierdistanz 2/9

#### Anwendungsbeispiele

Beispiel 1: Dubletten in Adressdatenbanken

Hein Blöd, 27568 Bremerhaven Hein Bloed, 27568 Bremerhafen Hein Doof, 27478 Cuxhaven

Beispiel 2: Produktsuche

Memori Stik

Beispiel 3: Websuche

eyjaföllajaküll semesta verien 2017



- Definition Editierdistanz ... alternativ: Levenshtein-Distanz
  - Gegeben zwei Zeichenketten x und y
  - ED(x, y) = die minimale Anzahl der folgenden Operationen, die man braucht, um x in y zu transformieren:

Einfügen (insert) eines Buchstabens

Ersetzen (replace) eines Buchstabens durch einen anderen

Löschen (delete) eines Buchstabens

 Die Position einer Operation ist die Stelle im String, an der etwas geändert wird ... siehe Beispiel nächste Folie

Positionen fangen hier und in der Folge mit **1** an, nicht 0 (weil das intuitiver ist bei der mathematischen Analyse)

## Editierdistanz 4/9

## DOOF SAUDOOF 5

#### Beispiel

$$-x = DOOF, y = BLOED ... ED(x, y) = ?$$

Danist Jahen nur enstmal nur gezeigt : ED(x,y) \le 4

## Editierdistanz 5/9



#### 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  $1 \le i \le j \le |x|$

Wie gesagt: Positionen / Indizes fangen mit 1 an, nicht 0

$$x = BLOED$$
,  $|x| = 5$ ,  $x[2..4] = LOE$ 

#### Editierdistanz 6/9



#### ■ Ein paar einfache Eigenschaften

```
- ED(x, y) = ED(y, x)
- ED(x, \varepsilon) = |x|
- ED(x, y) \ge abs(|x| - |y|) \qquad abs(x) = x > 0 ? x : -x
- ED(x, y) \le max(|x|, |y|)
- ED(x, y) \le ED(x[1..n-1], y[1..m-1]) + 1 \qquad n = |x|, m = |y|
```

Die Beweise sind alle einfache Zwei- oder Dreizeiler

Sehr gute Übung zur Klausurvorbereitung und zur Überwindung der Mathe-Phobie

## Editierdistanz 7/9



- Lösungsidee 1: möglichst viel "erhalten"
  - ED("VERIEN", "FERIEN") = △
     Einfach, weil die Zeichenketten größtenteils gleich
  - ED("SEMESTERFERIEN", "SEMESTERVERIEN") = 4Dito ... dann auch für längere Zeichenketten noch einfach
  - ED("MEXIKO", "AMERIKA") = 3
     Auch hier gibt es noch relativ große Übereinstimmung
  - ED("AAEBEAABEAREEAE", "RBEAAEEBAAAEBBAE") = \$\frac{2}{2}\$
     Spätestens hier wird es sehr schwierig mit dieser Idee

## Editierdistanz 8/9



- Lösungsidee 2: in zwei Hälften teilen
  - In zwei gleich große Hälften teilen und die dann jeweils rekursiv lösen

```
ED(GRAU, RAUM) = 2

ED(GR, RA) = 2

ED(AU, UM) = 2
```

 Keine gute Idee: die ED zwischen den Hälften hat nicht viel zu tun mit der ED zwischen den ursprünglichen Zeichenketten

```
Formal: wenn x = x_1x_2 und y = y_1y_2, dann ist im Allgemeinen nicht ED(x, y) = ED(x_1, y_1) + ED(x_2, y_2)
```

## Editierdistanz 9/9



- Lösungsidee 3: alternative rekursive Formel
  - Das Problem auf ein Problem "eins kleiner" zurückführen
  - ED("FERIEN", "VERIEN") = ED("FERIE", "VERIE") = 4
     Weil die Worte rechts mit dem selben Buchstaben aufhören
  - ED("SAUM", "RAUS") = ED("SAU", "RAU") + 1 = 2
     Weil eine optimale Folge ein replace am Ende macht
     Das gilt aber nicht immer, wenn sich die letzten beiden Buchstaben unterscheiden, zum Beispiel:
  - $2 = ED("RAUM", "GRAU") \neq ED("RAU", "GRA") + 1 = 3$ 
    - Mit einer etwas komplizierteren Formel kriegt man es aber hin, das sehen wir jetzt auf den nächsten Folien

## Rekursive Formel 1/7

INSCI,U) INSCI, A) INSCI, S)

DOOF -> UDOOF -> AUDOOF -> SAUDOOFU

mich manatan

INSCI,S) INSCI,A) INSCI,U)

DOOF -> SDOOF -> SAUDOOF

das int manatan

#### Monotonie

- Seien x und y unsere beiden Zeichenketten
- Seien  $\sigma_1$ , …,  $\sigma_k$  eine Folge von k = ED(x, y) Operationen, für  $x \to y$ , 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

– Eine Folge von Operationen heißt **monoton**, wenn die Position von  $\sigma_{i+1}$  ist ≥ die Position von  $\sigma_i$ , wobei = nur dann erlaubt ist, wenn beides "delete" Operationen sind

Eine Folge von delete Operationen mit **gleichen** Positionen braucht man zum Beispiel bei ED("saudoof", "doof")

```
SAUDOOF - AUDOOF - UDOOF - DOOF

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

## UNI FREIBURG

#### Hilfssatz zur Monotonie

- **Lemma:** Für beliebige x und y mit k = ED(x, y) gibt es eine monotone Folge von k Operationen für  $x \rightarrow y$
- Der Beweis ist Aufgabe 1 vom ÜB11
- Beweisidee 1: für den Fall k = 2 (zwei Operationen), muss man sich im Prinzip nur alle neun Kombinationen der drei Arten von Operationen anschauen
- Beweisidee 2: eine allgemeine nicht-monotone Folge lässt sich durch "Nachbartranspositionen" immer in eine monotone Folge derselben Länge überführen

Nachbartransposition = zwei Operationen hintereinander in nicht-monotoner Reihenfolge werden "umgedreht"

Rekursive Formel 3/7

FALL 1A: GRAU TRAU RAU RAU

TALL 15: BAOME TRAUME

FALL 1C: DOOF TRAUME

TALL 1C: DOOF T

– Wir betrachten die letzte Operation  $\sigma_k$ 

$$\sigma_1, \dots, \sigma_{k-1}: x \to z \quad \text{und} \quad \sigma_k: z \to y$$

$$Seien \quad n = |x| \quad \text{und} \quad m = |y| \quad \text{und} \quad m' = |z|$$

$$Man \quad beachte, \quad dass \quad m' \in \{m-1, m, m+1\}$$

- Fall 1:  $\sigma_k$  macht etwas "ganz am Ende" von z :

```
Fall 1a: \sigma_k = insert(m' + 1, y[m]) [dann ist m' = m - 1]
                                        [dann ist m' = m + 1]
Fall 1b: \sigma_k = delete(m')
Fall 1c: \sigma_k = \text{replace}(m', y[m]) [dann ist m' = m]
```

Wenn keines von den dreien der Fall ist, dann sind die letzten Zeichen von x und y (und z) gleich ... siehe nächste Folie



#### Fallunterscheidung

– Wir betrachten die letzte Operation  $\sigma_k$ 

```
\sigma_1, ..., \sigma_{k-1}: x \to z \quad \text{und } \sigma_k: z \to y
Seien n = |x| \text{ und } m = |y| \text{ und } m' = |z|
Man beachte, dass m' \in \{m-1, m, m+1\}
```

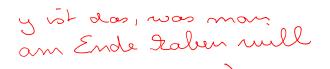
- Fall 2:  $\sigma_k$  macht nichts "ganz am Ende" von z

```
Dann z[m'] = y[m] und x[n] = z[m']

Damit \sigma_1, ..., \sigma_k: x[1..n-1] \rightarrow y[1..m-1] und x[n] = y[m]

\sigma_1 = \sigma_2 = \sigma_3 = \sigma_4
```

## Rekursive Formel 5/7



- Wir haben also einen dieser vier Fälle
  - Fall 1a (insert am Ende):  $\sigma_1$ , ...,  $\sigma_{k-1}$ : x[1..n] → y[1..m-1]
  - Fall 1b (delete am Ende):  $\sigma_1$ , ...,  $\sigma_{k-1}$ : x[1..n-1] → y[1..m]
  - − Fall 1c (replace am Ende):  $\sigma_1$ , ...,  $\sigma_{k-1}$ : x[1..n-1] → y[1..m-1]
  - Fall 2 (nichts am Ende):  $\sigma_1$ , ...,  $\sigma_k$ : x[1..n-1] → y[1..m-1]

Wichtig für die Formel auf der nächsten Seite: diese vier Fälle decken **alle** Möglichkeiten ab, andere gibt es nicht

## Rekursive Formel 6/7



- 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

      Jewels letste Operation
    - ED(x[1..n-1], y[1..m]) + 1
    - ED(x[1..n-1], y[1..m-1] +  $\frac{1}{1}$  ... falls x[n]  $\neq$  y[m]
    - ED(x[1..n-1], y[1..m-1] ... falls x[n] = y[m]

Jeder der vier führt zu einer möglichen Folge, und wir wissen, dass die minimale Folge dabei ist, deswegen das Minimum

- $F \ddot{u} |x| = 0 \text{ ist } ED(x, y) = |y|$
- $F\ddot{u}r |y| = 0 \text{ ist } ED(x, y) = |x|$

#### Alternative Sichtweise

Visualisieren einer Folge von Operationen, indem man
 x und y mit geeigneten "Lücken" untereinander schreibt

```
X DOO V M A N N I
Y B L O E D F R A U
```

insert = oben leer, darunter ein Buchstabe

delete = unten leer, darüber ein Buchstabe

replace = zwei ungleiche Buchstaben übereinander

- Wenn man da jetzt von links nach rechts durchgeht, hat man wieder genau eine monotone Folge
- Die rekursive Formel unterscheidet die vier Möglichkeiten in der letzten Spalte (insert, delete, replace, nix)

## Rekursive Implementierung 1/4



#### Code

- Mit der Formel von der Folie vorher können wir jetzt sehr leicht ein rekursives Programm schreiben
- Es funktioniert auch! (immerhin)

Aber es dauert unverhältnismäßig lange, selbst schon für relative kurze Zeichenketten

## Rekursive Implementierung 3/4



#### Laufzeit

Für die Laufzeit gilt folgende rekursive Formel

$$T(n, m) = T(n-1, m) + T(n, m-1) + T(n-1, m-1) + \Theta(1)$$

- Insbesondere:

$$T(n, n) \ge 3 \cdot T(n-1, n-1) \ge 3 \cdot 3 \cdot T(n-2, n-2) \ge ...$$

- Also:

$$T(n, n) \ge 3^{n-1} \cdot T(1, 1) \ge 3^{n-1}$$

Also **exponentielle** Laufzeit, wie auch schon bei der rekursiven Fibonacci-Berechnung

Insbesondere: Zeidentetten um 1 langer

- drei mal so lange

## Rekursive Implementierung 4/4

## UNI FREIBURG

#### Problem

Das rekursive Programm berechnet die gleichen ED
 Werte immer und immer wieder

Das hatten wir auch schon bei der Berechnung der
Fibonacci Zahlen gesehen (in Vorlesung 1b)

2.3. ED (ab, cd)

ED (s, cd)

ED (s, cd)

E | | | | | | | |

Lambdall

Bassifall

Lambdall

L

#### Literatur / Links

## UNI FREIBURG

#### Editierdistanz

In Wikipedia (Definition + Algorithmen)

http://en.wikipedia.org/wiki/Levenshtein distance

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