

Informatik II: Algorithmen und Datenstrukturen SS 2015

Vorlesung 11a, Dienstag, 7. Juli 2015
(Editierdistanz, Teil 1)

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, CC, Dijkstra)
- Offizielle Evaluation der Veranstaltung
- Wann ist ein Berg ein Berg

■ Inhalt

- Editierdistanz ... Motivation, Definition, rekursive Formel
 - Ü11, Aufgabe 1: Online-Evaluation (20 Punkte !!!)
- Zur Aufgabe 2 in der Vorlesung morgen mehr

Erfahrungen mit dem Ü10 1/4

■ Zusammenfassung / Auszüge Stand 7. Juli 12:00

- Aufgabe 1 (Laufzeit) gute Übung aber fiel vielen schwer
Siehe Erklärungen auf den nächsten drei Folien
- Aufgabe 2 (Dijkstra) spannend aber für viele auch schwieriger als erwartet (in den Details)
- Wie erwartet hatten einige Probleme mit BaWü
- Mein längster kürzester Pfad ist zu kurz
- Zu ... heiß ... zum ... Aufgaben Machen
- Abschließende Notizen auf dem Ü immer kryptischer, hoffentlich nicht klausur-elefant
- Rosskopf: Prominenz ca. 100m, Dominanz ca. 4km

Erfahrungen mit dem Ü10 2/4

und noch: Rekursion

■ Laufzeitbestimmung, Vorbetrachtung

- Grundsätzliches Vorgehen: die **Schleifen** anschauen

Nur da kann nicht-konstante Laufzeit verbraten werden

Dabei spielt es keine Rolle, ob es eine for Schleife oder eine while Schleife ist, wichtig ist, wie oft sie durchlaufen wird

- Achtung bei Funktionsaufrufen: innerhalb der Funktion können sich auch Schleifen befinden

Wenn man die Funktion nicht kennt, in der Doku nach Laufzeit schauen ... sonst plausible Annahmen machen

Zum Beispiel: beim Zugriff auf ein Element von einem Feld kann man davon ausgehen, dass es in $O(1)$ Zeit geht

■ Lösungsskizze Aufgabe 1

bei ungerichteten
Graphen: jede Kante
höchstens zweimal

- In unserem Code in `bfs.py` gibt es zwei Schleifen ... alle anderen Operationen haben konstante Laufzeit

Die äußere Schleife hat insgesamt konstante Kosten für jeden Knoten in jedem Level

Die innere Schleife hat insgesamt konstante Kosten für jede Kante, die von einem dieser Knoten ausgeht

- Durch das "visited" wird jeder Knoten und jede Kante nur **genau einmal** angeschaut
- Die Laufzeit ist also $\Theta(n' + m')$, wobei n' und m' die Anzahl Knoten und Kanten in der Zusammenhangskomponente sind, in der sich der Startknoten befindet

■ Laufzeitbestimmung, Nachbetrachtung

- Es gibt zwar Regeln (siehe Folie 4), aber die Laufzeitbestimmung ist kein einfacher schematischer Prozess

Vergleichbar mit einem mathematischen Beweis

In der Tat kann Laufzeitbestimmung beliebig komplex sein, selbst für einfache Algorithmen

Für BFS war es aber relativ einfach, und so wird es auch sein, wenn es in der Klausur dran kommt

Offizielle Evaluation der Veranstaltung

- Läuft über das zentrale **EvaSys** der Uni

- Sie sollten gestern (**Montag, 6. Juli**) eine Mail vom System bekommen haben

Falls nicht, bitte umgehend bei uns melden !

- Nehmen Sie sich bitte Zeit und füllen Sie den Bogen sorgfältig und gewissenhaft aus

Sie haben soviel Zeit in die Vorlesung investiert, dann können Sie auch 30 Minuten für die Evaluation aufwenden

Sie bekommen außerdem 20 wunderschöne Punkte dafür

- Uns interessieren besonders die Freitextkommentare

Wann ist ein Berg ein Berg

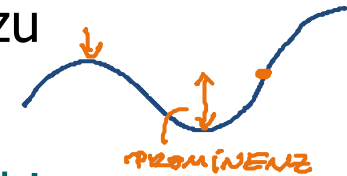
■ Dominanz und Prominenz

- Dominanz = Entfernung (Luftlinie) zum nächstgelegenen Punkt gleicher Höhe



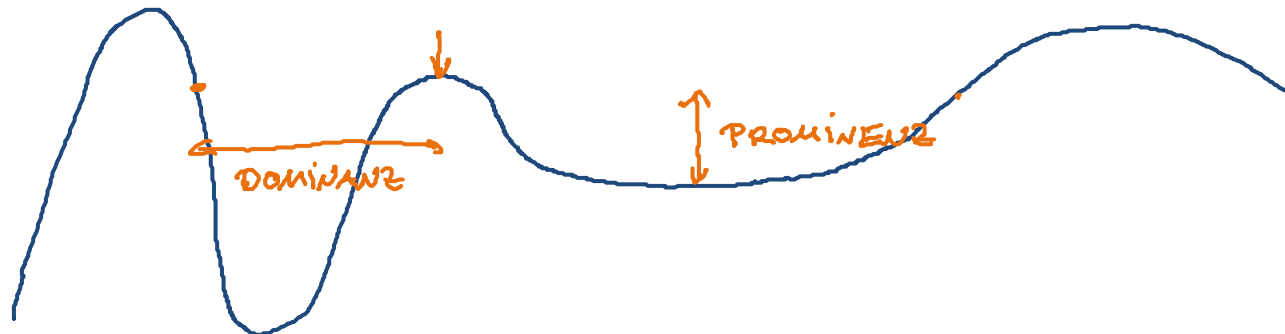
Intuitiv: wie weit der nächste höhere Berg weg ist

- Prominenz = minimaler Abstieg (Höhenmeter), um zu einem anderen Punkt gleicher Höhe zu kommen



Intuitiv: niedrigster Graben zu einem höheren Berg hin

- 1-Millionen-Euro Frage: kann dieser höhere Berg ein anderer sein für Dominanz und Prominenz ?



■ 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 der übernächsten Folie

Benannt nach dem Zeichenketten fressenden sogenannten **Edi-Tier**



■ 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

uni semesta verien 2015

- 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** eines Buchstabens (**insert**)
 - Ersetzen** eines Buchstabens durch einen anderen (**replace**)
 - Löschen** eines Buchstabens (**delete**)
 - Die **Position** einer Operation ist die Stelle im String, an der etwas geändert wird ... **siehe Beispiel nächste Folie**

Editierdistanz 4/16

■ Beispiel

DOOF \rightarrow BLOED

ED = ?

	0	1	2	3	...	
	D	O	O	F		
	B	O	O	F		REPLACE(0, B)
	B	L	O	F		REPLACE(1, L)
	B	L	O	E	F	INSERT(3, E)
	B	L	O	E	D	REPLACE(4, D)

Diese Folge
ist
MONOTON
(s. Folie 18)
 $0 < 1 < 3 < 4$

Das zeigt jetzt erstmal nur,
dass $ED \leq 4$

(Es könnte ja auch noch mit weniger gehen)

■ 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 $0 \leq i \leq j < |x|$

$$x = \overset{0}{B} \overset{1}{L} \overset{2}{O} \overset{3}{E} \overset{4}{D} \quad x[1..3] = LOE \quad |x| = 5$$

■ Ein paar einfache Eigenschaften

– $ED(x, y) = ED(y, x)$

– $ED(x, \varepsilon) = |x|$

– $ED(x, y) \geq \text{abs}(|x| - |y|)$ $\text{abs}(x) = x > 0 ? x : -x$

– $ED(x, y) \leq ED(x[1..n-1], y[1..m-1]) + 1$ $n = |x|, m = |y|$

$ED(x, y) = ED(y, x)$

Es gibt zu jeder Folge $x \rightarrow \dots \rightarrow y$ eine entsprechende Folge $y \rightarrow \dots \rightarrow x$ die dasselbe rückwärts macht

REPLACE \leftrightarrow REPLACE

INSERT \leftrightarrow DELETE

DELETE \leftrightarrow INSERT

■ Lösungsidee 1: möglichst viel "erhalten"

– $ED(\text{"VERIEN"}, \text{"FERIEN"}) = 1$

Einfach, weil die Zeichenketten größtenteils gleich

– $ED(\text{"SEMESTERFERIEN"}, \text{"SEMESTERVERIEN"}) = 1$

Dito ... dann auch für längere Zeichenketten noch einfach

– $ED(\text{"MEXIKO"}, \text{"AMERIKA"}) = 3$

Auch hier gibt es noch relativ große Übereinstimmung

– $ED(\text{"AAEBEAABEAREEEAEB"}, \text{"RBEAAEEBAAAEBBAEAE"}) = 2$

Spätestens hier wird es sehr schwierig mit dieser Idee

■ Lösungsidee 2: in zwei Hälften teilen

- In zwei gleich große Hälften teilen und die dann jeweils rekursiv lösen

$$ED(\underline{\text{GRAU}}, \underline{\text{RAUM}}) = 2$$

$$ED(\text{GR}, \text{RA}) = 2$$

$$ED(\text{AU}, \text{UM}) = 2$$

$$x = x_1 x_2 \quad y = y_1 y_2$$

Dann im Allgemeinen

$$ED(x, y) \neq ED(x_1, y_1) + ED(x_2, y_2)$$

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

■ Lösungsidee 3: alternative rekursive Formel

- Das Problem auf ein Problem "eins kleiner" zurückführen

$$ED("FERIEN", "VERIEN") = ED("FERIE", "VERIE")$$

Weil die Worte rechts mit dem selben Buchstaben aufhören

$$ED("SAUM", "RAUS") = ED("SAU", "RAU") + 1$$

Weil eine optimale Transformation ein replace am Ende macht

$$ED("RAUM", "GRAU") \neq ED("RAU", "GRA") + 1$$

Das gilt aber nicht immer, wenn sich die letzten beiden Buchstaben unterscheiden

Mit einer etwas komplizierteren Formel kriegt man es aber hin, das sehen wir jetzt auf den nächsten Folien

SAUDOOF →

→ BUDDOOF

■ Monotonie

- 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

- Eine Folge von Operationen heißt **monoton**, wenn die Position von σ_{i+1} ist \geq die Position von σ_i , wobei $=$ nur dann erlaubt ist, wenn beides **delete** Operationen sind

Eine Folge von delete Operationen mit gleichen Positionen braucht man zum Beispiel bei $ED(\text{"saudooof"}, \text{"doof"})$

$\overset{0}{S} \overset{1}{A} \overset{2}{U} \overset{3}{D} \overset{4}{O} \overset{5}{O} \overset{6}{F} \rightarrow \overset{0}{A} \overset{1}{U} \overset{2}{D} \overset{3}{O} \overset{4}{O} \overset{5}{F} \rightarrow \overset{0}{U} \overset{1}{D} \overset{2}{O} \overset{3}{O} \overset{4}{F} \rightarrow \overset{0}{D} \overset{1}{O} \overset{2}{O} \overset{3}{F}$
DELETE(0) DELETE(0) DELETE(0)

■ Hilfssatz zur Monotonie

- **Lemma:** Für beliebige x und y mit $k = ED(x, y)$ gibt es eine monotone Folge von k Operationen für $x \rightarrow y$
- Den Beweis lassen wir hier weg

Intuition: die Reihenfolge der Operationen ist im Grunde egal, also kann man sie auch monoton anordnen

Editierdistanz 12/16

Beispiele für die drei U-Fälle
Fall 1a: DOOF $\xrightarrow{\sigma_1 \sigma_3}$ BLOE $\xrightarrow{\sigma_4}$ BLOED

Fall 1b: TROTTEL $\xrightarrow{\sigma_1 \dots \sigma_5}$ IDIOTT $\xrightarrow{\sigma_6}$ IDIOT

Fall 1c: DEPP $\xrightarrow{\sigma_1, \sigma_2}$ NERP $\xrightarrow{\sigma_3}$ NERD

■ Fallunterscheidung

- Wir betrachten die letzte Operation σ_k

z.B. $\ell = 4$ DOOF BLOEF BLOEF BLOED
 $\sigma_1, \dots, \sigma_{k-1} : X \rightarrow Z$ und $\sigma_k : Z \rightarrow Y$

z.B. Seien $n = |X|$ und $m = |Y|$ und $m' = |Z|$

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

- **Fall 1:** σ_k macht etwas "ganz am Ende" von z :

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 $m' = m$]

Eine andere Möglichkeit als die drei gibt es nicht

Indizes fangen zwei mit 1 an und in der Folge

Beispiel zu Fall 2:
DOOF! → BLOED!
 $\sigma_1 - \sigma_4$

■ Fallunterscheidung

- Wir betrachten die letzte Operation σ_k

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

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

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

- **Fall 2:** σ_k macht nichts "ganz am Ende" von z

$$\text{Dann } z[m'] = y[m] \text{ und } x[n] = z[m']$$

$$\text{Damit } \sigma_1, \dots, \sigma_k : x[1..n-1] \rightarrow y[1..m-1] \text{ und } x[n] = y[m]$$

- 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

- $ED(x[1..n], y[1..m-1]) + 1$ *für das „insert“ am Ende*
- $ED(x[1..n-1], y[1..m]) + 1$ *für das „delete“ am Ende*
- $ED(x[1..n-1], y[1..m-1]) + \delta$ *für das „replace“ am Ende*

wobei $\delta = 1$ falls $x[n] \neq y[m]$, sonst $\delta = 0$

– Für $|x| = 0$ ist $ED(x, y) = |y|$

– Für $|y| = 0$ ist $ED(x, y) = |x|$

■ Rekursive Implementierung

- 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

- Editierdistanz

- In Wikipedia (Definition + Algorithmen)

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

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