

EINIGE BEMERKUNGEN ZU GLEICHHEITSTESTS UND DEREN KOMPLEXITÄT

Herrn Prof. Dr. St. SCHOTTLAENDER
in Dankbarkeit und Respekt zu seinem 60. Geburtstag

VON

W. LEX

In der - sit venia verbo - Kunst des Programmierens wird man des öfteren mit der Aufgabe konfrontiert, festzustellen, ob in einer Liste Elemente übereinstimmen; etwa, wenn ein durch seine Cayleytafel gegebenes Magma auf Kürzbarkeit hin zu untersuchen ist oder wenn in einer Datenbank, z. B. einer Personalkartei, gleiche Attributwerte bestimmt werden sollen. Im wesentlichen läuft das darauf hinaus, daß eine endliche Folge natürlicher Zahlen auf gleiche Glieder hin getestet wird.

Genauer - \mathbb{N} sei die Menge positiver ganzer Zahlen und $n = \{x \in \mathbb{N} \mid x \leq n\}$ für $n \in \mathbb{N}$ - legen wir fest

Definition 1: $\alpha = (a_x)_{x \in k}$ sei eine Folge auf n mit $k, n \in \mathbb{N} \setminus \{1\}$; ein Algorithmus γ heiße *Gleichheitstest* - für α -, falls γ die Wahrheit bzw. Falschheit von

$$(1) \quad \exists \kappa, \lambda \in k: \kappa < \lambda \wedge a_\kappa = a_\lambda$$

feststellt.

Für $k > n$ liefert das Schubfachprinzip sofort die Gültigkeit von (1). - Wir beschränken uns hier auf $k = n$, da der allgemeine Fall später ausführlicher dargestellt werden soll.

Nach einer vollständigen Analyse des üblichen Vergleichstests in I wird in II eine Klasse von Algorithmen betrachtet, die die Folge α nur einmal linear durchlaufen und bei minimalem λ in (1) abbrechen. - Gleichheitstests scheinen bislang in der Literatur nicht behandelt worden zu sein; für entsprechende Hinweise ist der Autor dankbar. ¹⁾

I

Der wohl nächstliegende Gleichheitstest für α mit $k = n$ besteht im sukzessiven Vergleich der einzelnen Folgenglieder mit den restlichen und Abbruch bei der ersten Gleichheit, genauer

Definition 2: Sei α wie in Definition 1 und $k = n$ ($n \in \mathbb{N} \setminus \{1\}$):

- α) ~~10~~ bezeichne den durch folgende PASCAL-Anweisungen gegebenen Algorithmus, wobei α im Feld A stehe, G eine Boolesche und I, J Integer-Variable sowie N eine Konstante für n seien:

```

G := false;
for I := 1 to N-1 do
  for J := I+1 to N do
    if A [I] = A [J] then  $\eta$ ;
1: if G then writeln ('Gleichheit')
   else writeln ('Alle verschieden');

```

mit η als Abkürzung für

```

begin G := true; goto 1 end .

```

¹⁾ Frl. cand. inf. M. Söding gilt mein aufrichtiger Dank für das Testen der Algorithmen und für kritische Durchsicht.

BEMERKUNGEN ZU GLEICHHEITSTESTS

β) V_n sei die *Durchschnittskomplexität* für 16, Gleichverteilung vorausgesetzt, also die mittlere Anzahl der zur Feststellung von Gleichheit oder Verschiedenheit (zur Erreichung der Marke 1 im Teil α)) notwendigen Vergleiche (Abfrage in Zeile 4 des obigen Programmausschnittes), wenn jede Folge α gleichwahrscheinlich ist.

Optimale Komplexität, 1, und pessimale, $\sum_{v=1}^{n-1} (n-v) = \binom{n}{2}$, liefern eine erste grobe Abschätzung

$$1 \leq V_n \leq \binom{n}{2};$$

genauer erhält man

Satz 1: Für $n \in \mathbb{N} \setminus 1$ gilt

$$(2) \quad V_n = n - R_n$$

mit

$$R_n \approx \frac{n!}{n^n} \sum_{v=0}^{n-1} \frac{v^v}{v!}$$

und

$$(3) \quad 0,581\ 976\dots = \frac{1}{e-1} < R_n \leq 1;$$

überdies erweist sich $\rho \approx (R_n)_{n \in \mathbb{N} \setminus 1}$ als streng monoton fallende, gegen $(e-1)^{-1}$ konvergierende Folge.

Beweis: Wir betrachten alle n^n n -tupel über n und bezeichnen mit A_n die Gesamtzahl der notwendigen Vergleiche; also

$$(4) \quad V_n = n^{-n} A_n.$$

Der Gleichverteilung wegen beschränken wir uns auf die etwa mit 1 beginnende Folgen, nennen die Gesamtzahl der notwendigen Vergleiche - bei $a_1 = 1$ - E_n und haben

$$(5) \quad A_n = n E_n.$$

Ferner sei B_n die Zahl der notwendigen Vergleiche bei

$$a_1 = 1 \wedge \exists v \in n \setminus 1: a_v = 1;$$

damit ergibt sich

$$(6) \quad E_n - B_n = (n-1) \cdot (n-1)^{n-1} + A_{n-1}.$$

b_v sei die Zahl der notwendigen Vergleiche im Falle $a_1 = 1 = a_{v+1}$ und $a_\mu \neq 1$ für $\mu \in v \setminus 1$ und $v \in n-1$; mit- hin gilt

$$b_v = v(n-1)^{v-1} \cdot n^{n-1-v},$$

unter Beachtung der sich durch Differentiation aus der Summe einer geometrischen Reihe ergebenden Summations- formel also

$$B_n = \sum_{v=1}^{n-1} b_v = n^n - (2n-1)(n-1)^{n-1}$$

und mit (4), (5) und (6) sowie $q_n \approx \frac{n-1}{n}$ für $n \in \mathbb{N}$ weiter

$$V_n = n + q_n^{n-1} (V_{n-1} - n),$$

was unschwer auf die Vermutung

$$V_n = n - \sum_{\mu=1}^n \prod_{v=\mu}^n q_v^{v-1}$$

führt, die sich mittels Induktion verifizieren läßt und nach einiger Umformung (2) impliziert.

BEMERKUNGEN ZU GLEICHHEITSTESTS

Durch Induktion beweist man

$$\sum_{v=0}^{n-1} \frac{v^v}{v!} \leq \frac{n^n}{n!} \quad (n \in \mathbb{N}),$$

was die obere Schranke in (3) liefert, und wegen²⁾

$$(7) \quad R_{n+1} = \left(1 + \frac{1}{n}\right)^{-n} (R_n + 1)$$

auch $R_{n+1} < R_n$, woraus die Konvergenz von ρ folgt. Mit $\lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n = e$ und (7) ergibt sich daher der Grenzwert von ρ zu $(e-1)^{-1}$.

II

Wie erwähnt behandeln wir nun einige Tests, die das zu testende n -tupel nur linear durchwandern. Um deren Komplexität zu beschreiben, geben wir

Definition 3: Mit α wie in Definition 1 und $k = n$ ($n \in \mathbb{N} \setminus \{1\}$) sei

$$m(\alpha) \triangleq \min \{ \lambda \in \mathbb{n} \mid \exists \kappa \in \mathbb{n} : \kappa < \lambda \wedge a_\kappa = a_\lambda \},$$

ferner

$$m_\alpha \triangleq \begin{cases} m(\alpha) \\ n \end{cases}, \text{ falls } \begin{cases} (1) \\ \rightarrow (1) \end{cases}$$

und

$$M_n \triangleq n^{-n} \sum_{\alpha \in \mathbb{n}^n} m_\alpha$$

(wobei $A^B \triangleq \{f \mid f : B \rightarrow A\}$).

Elementare Überlegungen und eine grobe Abschätzung liefern

²⁾Für den Hinweis, hier rekursiv vorzugehen, danke ich Herrn Dr. J. Schalmack herzlich [2].

Satz 2: Für alle $n \in \mathbb{N} \setminus 1$ gilt

$$M_n = (n-1)! \left(\sum_{v=1}^{n-2} \frac{v(v+1)}{n^v (n-v)!} + n^{3-n} \right)$$

und

$$(8) \quad 2 \leq M_n < \frac{n+3}{2}$$

(dabei sei $\sum_{v=1}^0 \dots \leq 0$).

In der nächsten Definition werden drei Algorithmen der oben genannten Art festgelegt: \mathcal{E} schöpft n aus, \mathcal{P} nutzt - ähnlich wie bei Gödelisierungen - die Eindeutigkeit der Primfaktorzerlegung, und \mathcal{Z}^3) gebraucht - wie auch \mathcal{P} - die Folgenglieder als Indizes.

Definition 4: α und $k=n$ seien wie in Definition 2, ferner werde der Programmteil von der Marke 1 an jeweils von dort übernommen ebenso wie \mathcal{Y} und die Bedeutung der einschlägigen Variablen und Konstanten, während der Typ der nicht eigens vereinbarten Variablen sinngemäß zu ergänzen ist. Durch die folgenden PASCAL-Anweisungen werden die erwähnten Verfahren beschrieben, und zwar

a) \mathcal{E} durch

```
G := false; Z := [1..N] - [A[1]];
for I := 2 to N do
  if A[I] in Z then Z := Z - [A[I]] else  $\mathcal{Y}$ ;
1: ...
```

b) \mathcal{P} durch

³⁾ Herrn Prof. Dr. J. Gillis danke ich aufrichtig für die freundliche Angabe eines Tests, der im wesentlichen mit \mathcal{Z} übereinstimmt [1].

BEMERKUNGEN ZU GLEICHHEITSTESTS

```
G := false; Q := Q div P[A[1]];
for I := 2 to N do
  begin D := P[A[I]];
    if Q mod D = 0 then Q := Q div D else  $\gamma$  end;
1: ... ,
```

wobei im Feld P die n ersten Primzahlen in natürlicher - oder auch beliebiger! - Reihenfolge und in Q deren Produkt stehen.

r) \mathcal{F} durch

```
G := false, W[A[1]] := false;
for I := 2 to N do
  if W[A[I]] then W[A[I]] := false else  $\gamma$ ;
1: ... ,
```

wobei jede Komponente des Feldes W zunächst auf "true" gesetzt ist.

Zum Abschluß vergleichen wir die eben definierten Tests mit dem in I behandelten Verfahren \mathcal{W} :

Satz 3: Bei einem Gleichheitstest γ für Folgen der Länge $k = n$ ($n \in \mathbb{N} \setminus \{1\}$) ergibt sich die mittlere Zeitkomplexität $K_n(\gamma)$ von γ unter Vernachlässigung der Aufbereitung der zusätzlich benutzten Felder und Mengen sowie eventueller additiver Konstanten zu

$$(9) \quad K_n(\mathcal{W}) = V_n v, \quad K_n(\gamma) = M_n c_\gamma$$

für $\gamma \in \{\mathcal{E}, \mathcal{P}, \mathcal{F}\}$ mit

$$v, c_\gamma \in \mathbb{R}^+ (\cong \{x \in \mathbb{R} \mid x > 0\});$$

ferner gilt

$$2v > c_y \Rightarrow \exists m \in \mathbb{N} \forall n \in \mathbb{N} (n > m \Rightarrow K_n(y) < K_n(\lambda))$$

für $y \in \{\mathcal{E}, \beta, \mathcal{I}\}$.

Beweis: (9) folgt unmittelbar aus der Definition der betreffenden Algorithmen.

Es sei $2v > c \approx c_y$ und m minimal aus \mathbb{N} mit

$$m \geq \frac{3c + 2v}{2v - c},$$

was für alle $n \in \mathbb{N}$ mit $n > m$ unter Beachtung von (9), (8), (2) und (3) sofort

$$K_n(y) = M_n c < \frac{n+3}{2} c < (n-1)v \leq V_n v = K_n(\lambda),$$

also die Behauptung liefert.

Numerisch empirische Bestätigung dieser Vergleiche sind geplant.

Referenzen

1. J. Gillis: Mündliche Mitteilung Mai 1987; Eremo dei SS. Pietro e Paolo, Bienno.
2. J. Schmalmack: Mündliche Mitteilung Herbst 1985, Clausthal.

W. Lex
 Institut für Informatik
 der Technischen Universität Clausthal
 Erzstr. 1
 D-3392 Clausthal-Zellerfeld