

Nochmals Gleichheitstests

von

Walter Klotz und Wilfried Lex

Abstract:

Procedures testing n -tuples for equal entries are called *equality tests*. The average case complexity is determined explicitly for various equality tests and asymptotic formulas are given.

A, B, C seien Mengen, dabei A geordnet und 0 enthaltend, und es sei

$$A^+ = \{a \in A \mid a > 0\}$$

und

$$C^B = \{f \mid f : B \rightarrow C\}$$

sowie $\mathbb{N} = \mathbb{Z}^+$ und

$$\underline{n} = \{x \in \mathbb{N} \mid x \leq n\} \quad (n \in \mathbb{N}).$$

In Programmen muß gelegentlich geprüft werden, ob eine Liste - etwa der Länge $n (\in \mathbb{N})$ - gleiche Einträge - etwa Zahlen aus \underline{m} mit $m \in \mathbb{N}$ - enthält [3]. Dies läuft auf die Frage nach der Injektivität von t aus \underline{m}^n hinaus; Algorithmen, die diese Frage beantworten, nennen wir *Gleichheitstests*.

Im folgenden werden einige Ergebnisse zur Komplexität solcher Tests ohne Beweise angegeben, die den Rahmen dieser Zusammenfassung sprengen würden und anderenorts dargestellt werden sollen [2]. - Da die expliziten Formeln für die Durchschnittskomplexität von Gleichheitstests i.a. schon relativ kompliziert sind und sich schwer miteinander vergleichen lassen, werden die Komplexitäten jeweils auch asymptotisch beschrieben.

Für den Trivialfall $m < n$ liefert das Schubfachprinzip sofort die Nichtinjektivität für jedes $t (\in \underline{m}^n)$; dennoch werden die Verfahren auch für $m < n$ untersucht.

I. Vergleichsalgorithmen

Von den insgesamt $\binom{n}{2}!$ Gleichheitstests, die auf direkten Vergleichen beruhen, betrachten wir hier lediglich zwei: \mathcal{V} arbeitet mittels „Vorwärtsvergleichen“, indem für $\nu = 1, 2, \dots, n-1$ jeweils $t(\nu)$ mit $t(\nu+1), t(\nu+2), \dots, t(n)$ verglichen und bei $t(\nu) = t(\nu+\rho)$ für $\rho \in \underline{n-\nu}$ abgebrochen wird, während \mathcal{W} „Rückwärtsvergleiche“ benutzt, also für $\nu = 2, 3, \dots, n$ jeweils $t(\nu)$ mit $t(\nu-1), t(\nu-2), \dots, t(1)$ vergleicht bis zum Abbruch bei $t(\nu) = t(\nu-\rho)$ ($\rho \in \underline{\nu-1}$).

Zur Definition der mittleren Komplexität von \mathcal{V} , i.e. Zahl der Vergleiche, sei $v_{m,n} : m^n \rightarrow \mathbb{N}_0$ für $m, n \in \mathbb{N}$ sowie für $n > 1$ weiter $v_{m,n}(t)$ die Zahl der notwendigen Vergleiche, um mittels \mathcal{V} das kleinste j aus \underline{n} mit

$$\exists k \in \underline{n} \setminus j : t(j) = t(k)$$

zu finden, falls ein solches j existiert, und $v_{m,n}(t) = \binom{n}{2}$ andernfalls; schließlich sei $v_{m,n} = (t \mapsto 0)$. Damit ist

$$V_{m,n} = m^{-n} \sum_{t \in m^n} v_{m,n}(t) \quad (m, n \in \mathbb{N})$$

die *Durchschnittskomplexität* von \mathcal{V} .

In deutlicher Verallgemeinerung von [3], Satz 1, S. 143, erhält man

Satz 1:

Für $m, n \in \mathbb{N}$ gilt

a) $V_{m,n} = m - R_{m,n}$ mit

$$R_{m,n} = \frac{m!}{m^n} \cdot \begin{cases} \left(\frac{m-n}{(m-n)!} + \sum_{\nu=1}^n \frac{(m-\nu)^{n-\nu}}{(m-\nu)!} \right) \\ \sum_{\mu=1}^{m-1} \frac{(m-\mu)^{n-\mu}}{(m-\mu)!} \end{cases} \quad \text{für } m \begin{cases} \geq \\ < \end{cases} n.$$

b) für $R_n = n - V_{n,n}$: $R_n = \frac{n!}{n^n} \sum_{\nu=0}^{n-1} \frac{\nu^\nu}{\nu!}$,

$$1 = R_1 = R_2 > R_3 > \dots > \frac{1}{e-1} \text{ ,}$$

$$\lim_{n \rightarrow \infty} R_n = \frac{1}{e-1} = 0,581\,976\dots$$

c) I. $q \in]0, 2[\wedge m \leq n^q \Rightarrow V_{m,n} - m \in O\left(\frac{m}{n}\right)$,

II. $c \in \mathbb{R}^+ \wedge m - cn^2 \in O(n) \Rightarrow V_{m,n} - c(1 - e^{-\frac{1}{2c}})n^2 \in O(n)$,

III. $q \in]2, \infty[\wedge m \geq n^q \Rightarrow V_{m,n} - \binom{n}{2} \in O\left(\frac{n^4}{m}\right)$.

$w_{m,n}$ und $W_{m,n}$ seien für \mathcal{W} ganz analog zu $v_{m,n}$ und $V_{m,n}$ für \mathcal{V} definiert. Bereits bei $m = n = 4$ differieren die mittleren Komplexitäten von \mathcal{V} und \mathcal{W} , und zwar um mehr als 1,46%:

$$W_{4,4} = \frac{101}{32} = 3,15625 < 3,203125 = \frac{205}{64} = V_{4,4}.$$

Allgemeineres über diese Abweichung liefert zusammen mit Satz 1 der

Satz 2:

Für $m, n \in \mathbb{N}$ gilt:

$$a) W_{m,n} = m^{-n} \left(\frac{1}{2m} \sum_{\nu=1}^n \nu(\nu^2 + 1) \binom{m}{\nu} \nu! m^{n-\nu} + \binom{m}{n} n! \left(\binom{n}{2} - \frac{n}{2m} (n^2 + 1) \right) \right),$$

für $m < n$ speziell

$$W_{m,n} = m + 1 - \frac{m!}{2m^m} \sum_{\mu=0}^m \frac{m^\mu}{\mu!}.$$

- b) I. $q \in]0, 2[\wedge m \leq n^q \Rightarrow W_{m,n} - m + \frac{1}{2} \sqrt{m \frac{\pi}{2}} \in o((\ln m)^2),$
 II. $c \in \mathbb{R}^+ \wedge m - cn^2 \in o(n) \Rightarrow W_{m,n} - c(1 - e^{-\frac{1}{2c}})n^2 \in o(n),$
 III. $q \in]2, \infty[\wedge m \geq n^q \Rightarrow W_{m,n} - \binom{n}{2} \in o\left(\frac{n^4}{m}\right).$

Die Beweise der in den Sätzen 1 und 2 vorgestellten Resultate sind l anglich; sie benutzen u.a. Rekursionsformeln f ur Binomialsummen der Form

$$\sum_{k=0}^n k^p \binom{m}{k} k! m^{n-k} \quad (n, p \in \mathbb{N}_0, m \in \mathbb{R}).$$

und deren asymptotisches Verhalten [1].

Die Frage nach dem „besten“ und „schlechtesten“ Vergleichsalgorithmus, d.h. nach dem mit der kleinsten und gr oten Durchschnittskomplexit at, ist offen; ebenso steht eine Klassifizierung aller Vergleichstests noch aus.

II. Lineare Tests

Nat urlich gibt es auch Gleichheitstests, bei denen das zu testende n -tupel t nur einmal linear durchwandert werden mu, etwa indem man die Glieder von t als Indizes f ur ein Boolesches Feld B der L ange m benutzt, dessen Komponenten zun achst alle auf „false“ gesetzt sind [3], Def. 38, S. 147; f ur $\nu = 1, 2, \dots, n$ wird nun $B[t(\nu)]$ auf „true“ gesetzt, es sei denn $B[t(\nu)] = true$, in welchem Falle abgebrochen wird. Es handelt sich also i. w. um ein „Hashen ohne Speichern“. – Solche Tests m ogen *linear* heien.

Um ihre Komplexit at zu beschreiben, sei $l_{m,n} : \underline{m}^n \rightarrow \mathbb{N}$ mit

$$l_{m,n}(t) = \begin{cases} n + 1, & \text{falls } t \text{ injektiv ist} \\ \min\{\nu \in \underline{n} \mid \exists \mu \in \underline{\nu - 1} : t(\mu) = t(\nu)\}, & \text{sonst} \end{cases}$$

f ur $t \in \underline{m}^n$; dann ist

$$L_{m,n} = m^{-n} \sum_{t \in \underline{m}^n} l_{m,n}(t)$$

die dazugeh orige *Durchschnittskomplexit at*.

Mit ähnlichen Methoden wie bei den Sätzen 1 und 2 läßt sich zeigen

Satz 3:

Für $m, n \in \mathbb{N}$ gilt:

$$a) L_{m,n} = m! \cdot \begin{cases} m^{-m} \sum_{\mu=0}^m \frac{m^\mu}{\mu!} \\ m^{-n} \sum_{\nu=0}^n \frac{m^{n-\nu}}{(m-\nu)!} \end{cases} \quad \text{für } m \begin{cases} \leq \\ > \end{cases} n.$$

$$b) \text{ I. } m \ln m \leq n^2 \quad \Rightarrow L_{m,n} - \sqrt{m \frac{\pi}{2}} \in o((\ln m)^2),$$

$$\text{ II. } c \in \mathbb{R}^+ \wedge m - cn^2 \in o(n) \quad \Rightarrow L_{m,n} - n \int_0^1 e^{-\frac{x^2}{2c}} dx \in o(1),$$

$$\text{ III. } q \in]2, \infty[\wedge m \geq n^q \quad \Rightarrow L_{m,n} - (n+1) \in o\left(\frac{n^3}{m}\right).$$

Hierbei ist die Zeit zur Anlage des Feldes und der zeitliche Aufwand zum Aufsuchen der einzelnen Komponenten noch nicht berücksichtigt, so daß die Frage bleibt, welches Verfahren „in praxi“ das beste sei, zumal ja auch noch andere Tests existieren.

Literatur

1. W. Klotz: On Certain Binomial Sums. Erscheint in:
„Festschrift“ in honour of Gerhard Ringel;
ed. R. Bodendiek, R. Henn. Physica-Verlag, Heidelberg 1990.
2. W. Klotz, W. Lex: Testing n -tuples for Equal Entries.
Eingereicht bei „Discrete Appl. Math.“.
3. W. Lex: Einige Bemerkungen zu Gleichheitstests und deren Komplexität.
Publication de l'Institut de Recherche Mathématique Avancée,
348/S-17 (1988). Actes du Séminaire Lotharingien de Combinatoire,
17^e Session: 27 - 30 mai 1987, Eremo SS. Pietro e Paolo,
Bienna. S. 141 - 148.

Walter Klotz
Institut für Mathematik

Wilfried Lex
Institut für Informatik

der Technischen Universität Clausthal
Erzstr. 1
D-3392 Clausthal-Zellerfeld