

6 Q-7

文書検索向き数値検索方式の提案

志村隆則 川口久光 加藤寛次 畠山 敦 秋沢 充
(株) 日立製作所 中央研究所

1. はじめに

文書や書類を検索する場合、その文書中の数値データに関して大小比較をしたり、範囲を指定して検索したいという要求がある。一般に、文章内の数値データは左詰めになっているので、数値列の1文字目の数値を読み出した段階では、その数値列が何桁の数値かわからない。最終桁の次の文字、すなわち数値列に続く非数値文字を検出して初めて、その数値文字列が何桁かが分かる。したがって、高速な数値検索機能を実現するためには、文字列を1文字ずつ読み出しながら逐次大小比較処理を行い、数値文字列を読み終わった段階で、直ちに比較判定結果を出力できる数値検索方式が必要となる。

2. 数値条件分割方式の提案

数値検索機能を実現する方法として、オートマトンを用いた方式がある。この方式では、検索したい範囲(以下、数値条件という)の数値全てに対して状態遷移パスを割り付ける必要がある。そのため、状態遷移テーブルが大規模になり、オートマトン生成時間が長くなるという問題がある。

本報告では、数値条件を桁数ごとに分割して、各桁ごとに数値検索を行なう数値条件分割方式を提案する。本方式によれば、オートマトンが簡単になるので、状態遷移テーブルが小さくなり、オートマトンの生成時間も短縮することが可能となる。

数値条件の分割方法は、以下の通りである。

(1) 数値条件の下限値と上限値の桁数が2桁以上異なる場合

数値の範囲を次の3つの範囲に分割する。

- ① 下限値から下限値と同じ桁数の最大数値まで。
- ② 下限値より1桁大きい数値の最小値から上限値より1桁小さい数値の最大値まで。
- ③ 上限値と同じ桁数の数値の最小値から上限値まで。

(2) 数値条件の下限値と上限値の桁数が1桁異なる場合

- (1) の①と③の2つの数値範囲に分割する。
- (3) 数値条件の下限値と上限値が同じ桁数の場合数値の分割は行わない。

以上のように、数値条件を分割することにより、分割された各数値条件は特定の桁数の数値になるため、それぞれのオートマトンが簡単になる。

具体的に、

$$12 \leq K \leq 345 \quad (1)$$

の数値条件を例に取り説明する。この例では、数値条件の下限値と上限値の桁数が1桁異なるので、(2)の分割方法に従い、次の2つに分割する。

$$12 \leq K1 \leq 99 \quad (2)$$

$$100 \leq K2 \leq 345 \quad (3)$$

オートマトンは、図1のように極めて簡単になるため、オートマトン生成時間が短くなる。このオートマトンは、分割条件ごとにマイクロプロセッサを用いて並列に処理することによって、高速化が実現できる。

3. おわりに

数値条件分割方式を提案した。この方式の特徴は、以下の通りである。

- (1) オートマトンが簡単になり、オートマトン生成時間を短くできる。
- (2) 分割した数値条件ごとに並列に照合処理することにより、高速化が実現できる。

【参考文献】

[1] L.A. Hollar, "Text Retrieval Computers," Computer, Vol. 12, No. 3, Mar. 1979, pp. 40-50.
[2] A.V. Aho and M.J. Corasick, "Efficient String Matching," Communications of the ACM, Vol. 18, No. 6, June 1975, pp. 333-340.

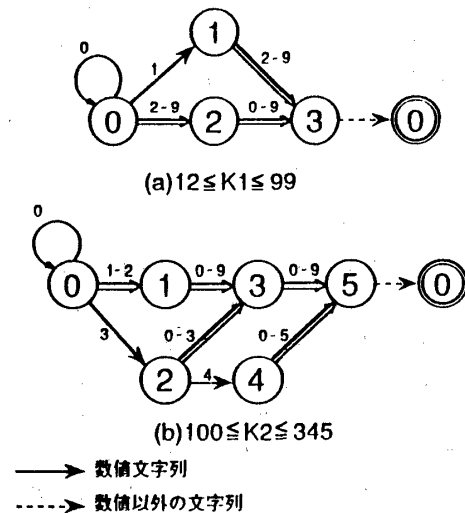


図1. 数値条件分割方式の状態遷移図($12 \leq K2 \leq 345$)