

区分木ソート

仁尾 都

明星大学経済学部

nio@mi.meisei-u.ac.jp

計算量が $O(n \log n)$ で、作業メモリが $\Theta(n)$ であり、かつ平均 `cpu` 時間がクイックソートと同程度のソートアルゴリズムを提案する。本ソートは逐次追加されるキーを平衡二分探索木の内点が所持する閾値によって定まるキーカバー範囲(区分)を持つ葉に割り当てる。葉に割り当てられたキーの個数が一定数より大きくなると、内点を新しく設け、その閾値を決定し、葉を二分し、二分木を平衡化する。葉のキーは挿入ソートとマージソートでソートされる。

Interval Tree Sort

Misato Nio

Faculty of Economics, Meisei University

nio@mi.meisei-u.ac.jp

Interval Tree Sort Algorithm, whose computing complexity $O(n \log n)$ is superior, and memory size $\Theta(n)$ inferior, but average sort time nearly equal to quick sort, is proposed. Every key additively given to the algorithm is assigned to one of leaves of a balanced binary search tree. When the number of keys assigned to a leaf exceeds the fixed number, the leaf is divided into two so that the tree is rebalanced. Keys in a same leaf are sorted by insertion sort and merge sort.

1. はじめに

殆どの計算機アルゴリズムの教科書で最初に出てくるテーマはキー群をその値の大小に従って整列するソートアルゴリズムである。かようにソートは最もポピュラーなアルゴリズム研究分野となっており、その研究開発の歴史は長い。ソートは多くの研究分野だけでなく、会計処理や統計処理の事務処理分野などの日常的業務にも繰り返し活用されており、常により高速なソート技法が望まれている。現在平均的な処理能力が最も高いと一般に認められているソートアルゴリズムはクイックソートである。しかしながら欠点としてその性能がキーの入力順序に大きく依存し、最悪計算量は $O(n^2)$ 、スタック領域も $O(n)$ となることに起因する不都合な事態が発生するなど、安心した利用ができない。したがって、クイックソートより平均処理速度が高速で、しかも計算量が $O(n \log n)$ で、作業メモリも小さく、かつ常に安心して用いることのできるソートアルゴリズムが求められている。本報告では、作業メモリを最大 $(2n + 6n/k)$ だけ必要 (k は $\log n$ 程度)

とするものの、計算複雑度が $O(n \log n)$ であり、かつ平均 **cpu** 時間がクイックソートと同程度で、ソートを安全確実に実行できるアルゴリズムを提案する。

2. 従来の研究

ソートアルゴリズムの研究分野には今までに非常に多くの努力とアイデアが注ぎ込まれてきた。その結果、多種多様のアルゴリズムが生まれた。しかしながら、他のどのアルゴリズムも、その平均 **cpu** 時間でクイックソートを凌駕できない。そのため、クイックソートの欠点である $O(n^2)$ のケースに対処するために、クイックソートの改良が試みられている。主な改良ポイントはピボット選択法の改良であるが、それだけでは $O(n^2)$ は改良されない。ピボット選択法として平均的に優れているとされる **Median-of-3** も **universal** な分布に対して $O(n^2)$ となり、その出現頻度も低くないことを指摘した **D.R.Musser** はクイックソートとヒープソートを使い分ける **Introsort** [1] を提案し、その性能を最悪でもヒープソート並みに抑えることができるとしている。また、平均処理速度がクイックソートよりも優れているとする河村らの多重分割ソート [2] があるが、これも内部にはクイックソートが用いられている。多重分割ソートはバケットソートに分類できる。ソートの前処理段階で n に応じてバケットの数を決定し、バケットの間の閾値をキーからランダムサンプリングすることにより、バケットに振り分けられるキーの平均個数を抑制することを狙ったものであり、平均処理性能をクイックソートより向上させることに成功している。しかしながら、ひとつのバケットに残りのキーが入ってしまうケースがあり、それらのキーをクイックソートでソートしているので依然として $O(n^2)$ となる。このように、クイックソートを凌駕するソートアルゴリズムの開発もまたクイックソートをベースとすることが行われてきた。

クイックソートをベースにしていないソート法として今回提案する区分木ソートは構造化された挿入ソートであると見做すことができる。キーをひとつずつ取り出してはソート済みのキー列 A に挿入ソートで挿入を行い、 A を徐々に長くし長大なキー列を作成する。しかるに挿入ソートは $O(n^2)$ であり、そのままでは使えない。そこで A を容量が一定 k の複数区分に分ける。これによって局所的に挿入ソートは $O(k^2)$ に抑制できる。区分を平衡二分探索木の葉に対応させて管理することによって、各キーがどの区分に属するかの処理は $O(\log(n/k))$ になる。葉内のキー数が多くなると葉を二分する。最後に二分木を辿って A を取り出す。

すでに二分木ソート (**Binary tree sort**) と呼ばれるソートアルゴリズムが公知である。このアルゴリズムは逐次追加的に与えられたキーを二分探索木に挿入し、最後に二分木を辿ってキーを整列させる。本来、二分木ソートは $O(n^2)$ であるが、キー毎に木を常に平衡に保つ処理を追加すれば $O(n \log n)$ となる。しかしながら、この処理の負担が大きすぎると考えられてきており、ヒープソートやマージソートなどの他の $O(n \log n)$ アルゴリズムと同列に比較されることが殆どなかった。二分木ソートではひとつのキーを

ひとつの内点や葉には一つのキーだけしか保持させないため、木が深くなる。そのため木の平衡化処理の回数がキーの数だけ必要となり、負荷が大きくなる。区分木ソートでは互いに近く値を持ったキーをひとつの区分として捕らえ、区分単位で葉に記憶し、また区分内のキーの数が区分の容量 $2k$ 以上になったときに葉を2分し、このときにのみ木の平衡化処理をすることとした。これにより、2分木の深さは浅くなり、平衡化処理回数は削減できる。

区分木ソートを **Library sort** [3] と対比できる。この手法では、挿入ソート列に所々に隙間を持たせ、新しく挿入するキーの格納スペースとする。挿入位置の決定方法は2分法による。これにより、高い確率的で $O(n \log n)$ を実現する。この点、区分木ソートは葉に固定長 $2k$ のスペースを設けており、同様に隙間は常に確保される。

区分木ソートの葉をバケットとして解釈すると、上記に述べた多重分割ソート法と比較することができる。多重分割ソート法では、ソート開始時にバケットの数とそれらの閾値を決定し、ソートが完了するまで固定する。区分木ソートではあくまでキーは逐次添加式で受け取り、区分の細分化は最後まで続けるので区分の閾値は常に変動する。

従来から2分木ソートで木を動的に平衡化し、またはバケット法でバケットを動的に平衡的に細分化することは大きな処理負荷になると考えられてきた [4]。平衡化の処理負荷を軽減する工夫を施した区分木ソートは、最悪計算量が $O(n \log n)$ であることと、平均 **cpu** 性能がクイックソートと同等であるという2つの長所を得た。

2. 区分木ソートアルゴリズム

(1) 区分木のデータ構造

区分木は平衡2分探索木であり、入力キーの領域とは別の作業メモリの中に作成される。木の根と各内点には2分探索木と同様、閾値を持たせ、これにより葉の担当区分範囲を定める。したがって、内点 p の左側に位置するすべての葉の担当区分範囲の最大値は p の所持する閾値よりも小さい。すべての葉は固定の長さ $2k+1$ のバッファを持っており、更にそれは長さ k の前半バッファと長さ k の後半バッファに分かれている。ソート途中の定常状態では、すべての葉の前半バッファには常にソート済みの k 個のキーが満杯状態で格納されており、後半バッファには高々 k 個のソート済みのキーが格納されている状態になっている。中間の $k+1$ 番地には後半の格納キー数が記憶される。図1は $k=3$ とした時の昇順ソートのための区分木の状態例で、11から42まで連続する整数のキーが処理され、木の内点と葉に蓄えられた有様を示す。各葉の前半と後半のキーは最終段階でマージされて出力される。キーが逐次添加されるにしたがって次第に葉のキーをカバーする範囲は細分化され、葉の数は増え、区分木の深さも深くなる。例えば19の閾値をもった内点の右の葉にもう1個のキーが追加されると葉は分割される。

(2) 区分木ソートのアルゴリズム

$key(i) (1 \leq i \leq n)$ を昇順ソートするとする。葉の容量 $2k+1$ の k は外部指定

である。

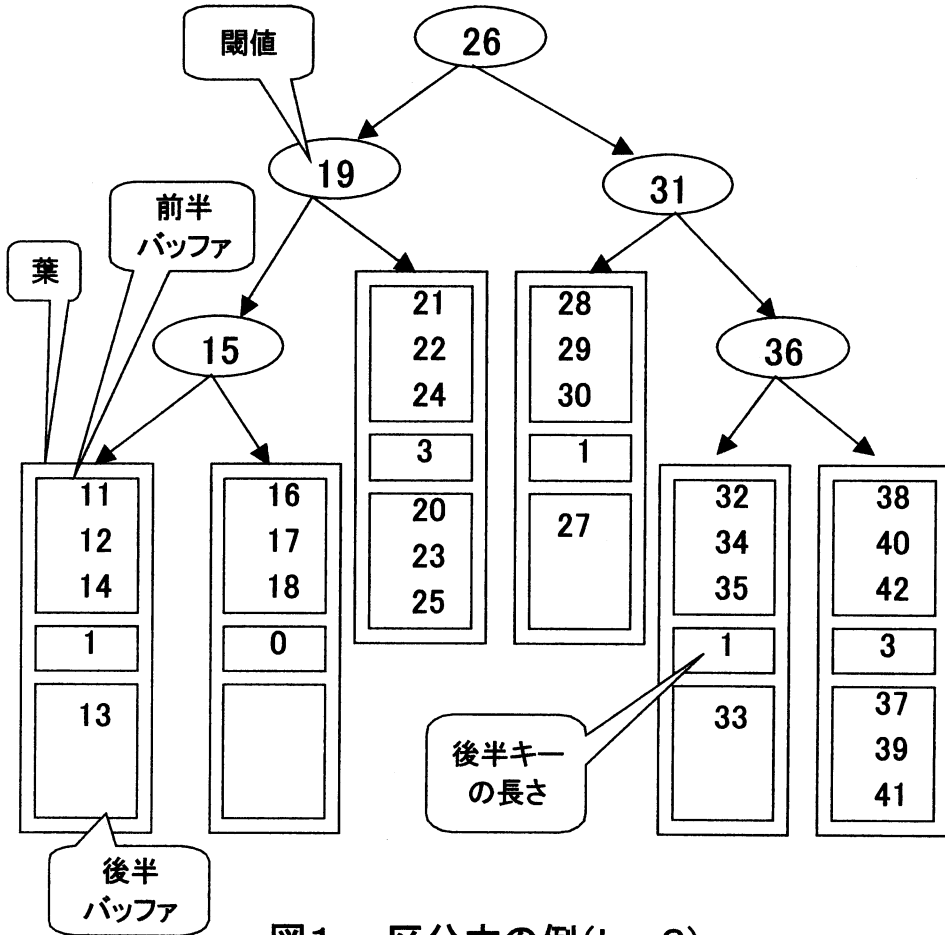


図1. 区分木の例($k=3$)

[手順1] 最初の $2k+1$ 個のキーを何らかの方法 (例: 挿入ソート) でソートする。もし $n > 2k+1$ なら、根を作成し、前半 k 個を左の葉の前半に、 $k+1$ 番目を根の閾値にし、最後の k 個を右の葉の前半に格納する。他の場合は終了。

[手順2] for $i=2k+2$ to n

[手順2. 1] 2分探索木を辿りキー ($key(i)$) が所属する葉 (i) を決定。

[手順2. 2]

if (葉 (i) の後半の格納キー数 $< k$) ([手順 i] キーを後半に挿入ソート)

elseif (葉 (i) の後半の格納キー数 $= k$) ‘葉 (i) の後半が満杯のとき’

([手順 ii] 内点 t とひとつの葉 a を新しく作成する。キーを葉 (i) の後半に挿入ソートし、結果を葉 (i) の前半とマージし、その結果の前半 k 個を葉 (i) の前半に、 $k+1$ 番目を t の閾値にし、最後の k 個

を a の前半に格納する。 t を葉 (i) の親の子とし、 t の左の葉を葉 (i) とし、 t の右の葉を a とする。

〔手順iii〕 2分木を平衡化する)

〔手順3〕 2分探索木に沿って順次葉のキーと内点の閾値を取り出し、 $key(i)$ の格納エリアにソート結果を上書きする。 この際、各葉の前半のキーと後半のキーはマージして出力する。

3. 計算複雑度

以下、ひとつのキー当たりの最悪計算量とその繰り返し回数を述べる。

〔手順1〕 $=O((2k+1)^2)$ 。理由： $2k+1$ 個のソートは挿入ソートと仮定。

——繰り返し回数は全部で1回。

〔手順2. 1〕 $=O(\log(n/k))$ 理由：葉の数は $O(n/k)$ 。平衡化は2色木を採用し、その深さは $O(2\log(n/k))$ 。——全部で $n-2k-1$ 回。

〔手順i〕 $=O(k)$ 。——全部で $(n-2k-1 - ((n-2k-1)/k))$ 回。

〔手順ii〕 $=O(k)$ 。理由：後半は k 個のキーに対し1個のキーを挿入ソート ($O(k)$)。その結果を前半とマージソート $O(2k+1)$ 。その結果を葉 (i) の前半と a の前半に格納 ($O(2k+1)$)。

——全部で $(n-2k-1)/k$ 回。

〔手順iii〕 $=O(\log(n/k))$ 理由：木は2色木で平衡化すると仮定。

〔手順3〕 $=O((n/k) + n)$ 理由：すべての内点 ($O(n/k)$ 個) は3回スキャンされる。各葉はその後半がすべて空である場合が最悪なので、手順3での前半と後半のバッファ間のマージのための計算量は考慮する必要はない。すべての葉の $O(n)$ のキーが出力バッファに移動する。——全部で1回。

したがってすべての合計のオーダーは $O(n \log n + nk)$ となり、 k の値を $\log n$ と同程度にすることにより、 $O(n \log n)$ とすることができる。

作業メモリは、1つの内点のために長さ5の管理メモリが必要である ($O(5n/k)$)。またすべての葉のバッファの合計が $O(2n + n/k)$ であるので、作業メモリは $O(2n + 6n/k)$ となる。

4. 計算機実験

区分木ソートの性能評価のために、区分木ソート、クイックソート及びヒープソートを計算機に実装した。区分木ソートのインターフェースを $its(key,n)$ とし、出力先も key とし、作業メモリはルーチン its の中で n に応じて動的に確保する。クイックソートは $quicksort(key,from,to)$ とし、コーディングを付録に掲載する。ピボットは3つのキーの平均とし、再帰実行時にソート対象数が12個以下になった場合は挿入ソートを行う。12

とした理由はその場合の処理速度が最も高速であったためである。ヒープソートは `quicksort(key,n)` としたが、そのコーディングは紙面数の都合により掲載できない。共に、使用した言語は EXCEL VBA で、計算機の `cpu` は PentiumIV (2 GHz) である。表 1 に各種条件下での、区分木ソートの `cpu` 時間の対クイックソート比と対ヒープソート比を掲載した。ただし、 $k = 1.2$ とした。従って比が 1 より小さい場合、2 分木ソートのほうが優れていることを示す。キーはすべて実数である。一様乱数 ($0 < \text{乱数} < 1$) は VBA の乱数発生関数 `Rnd` を用いた。標準偏差が与えられた正規分布の乱数はボックスミュラー法で作成した。

区分木ソートにとって都合の最も悪いケースは、キー出現順序が殆ど降順か昇順である場合である。なぜなら、そのような時は最近に作られた葉にキーが連続的に配分され、その葉が更に分割され、木の平衡化作業が連続的に繰り返され、他の葉の後半は空白のまま残されてしまうからである。またヒープソートにとってもキーが昇順や降順の場合はやはり最悪のケースになる。2 分木に対しキーを追加し、あるいは削除する場合、常に根から葉に至るまで、すなわち木の深さと同じ回数だけ木を上下するという探索処理を強いられるためである。そこで、昇順と降順のケースを取り上げた。昇順のキーはまず指定された分布のキーを乱数で作成し、次にそれをソートすることによって得た。降順のキーも同様である。

表 1. 区分木ソートの `cpu` 時間比 (対クイックソート, 対ヒープソート) ($N(m, \sigma)$: 正規分布)

n	キーの分布	一様分布	$N(0, 10^{-3})$	$N(0, 10^3)$	$N(0, 10^{-3})$
	キーの入力順序	ランダム	ランダム	昇順	降順
10^3	対クイックソート比	0.95	0.98	1.73	1.75
	対ヒープソート比	0.49	0.52	0.36	0.56
10^4	対クイックソート比	0.90	0.90	1.80	1.82
	対ヒープソート比	0.47	0.46	0.34	0.55
10^5	対クイックソート比	0.85	0.85	1.85	1.86
	対ヒープソート比	0.43	0.43	0.33	0.55
10^6	対クイックソート比	0.86	0.85	1.91	1.91
	対ヒープソート比	0.41	0.41	0.32	0.54

5. 結果の検討

表 1 を概観すると、キーの入力順序がランダムの場合、区分木ソートの計算機による実際の平均 `cpu` 時間はクイックソートより短く、また、区分木ソートはヒープソートに比べ 0.56 以下となっている。n が増加するにしたがってその比が向上する理由は、n が小さい間はクイックソートの再帰的 `Call` の回数が少なくその処理負荷の割合が軽いためである。

キー出現順序が降順か昇順であるような区分木ソートにとって最悪のケースでは、クイックソート比で最大1.9倍程度となりかなり劣る。しかし表には記載していないが、区分木ソートでは昇順や降順のキーの場合と一様分布の場合のcpu時間の比は1.24倍に収まっており、本手法が $O(n \log n)$ であることの利点が現れている。クイックソート比では最大1.9倍にもなる理由は、クイックソートで採用したピボットを最初と中間と最後のキーの平均する方法 (Median-of-3) を採用したことにより、昇順・降順キーに対する性能はランダム順のキーの場合に比べ、逆に大幅に性能が向上しているためである。しかしながらピボットの決定方法が如何に定められようとも、その決定方法にもっとも都合の悪い入力データが存在する。Median-of-3法に対して $O(n^2)$ となるようなMedian-of-3 killer [1] のケースを今回は取り上げていない。

平均cpu時間の実験結果は使用言語や計算機環境, 双方のコーディング方法に依存し, 特に今回クイックソートの再帰的 call の負荷軽減対策として挿入ソートを併用するだけでは不十分の可能性がある, 区分木ソートの平均cpu時間がクイックソートより優れているとは断定できない。この程度の差であれば区分木ソートの平均cpu時間は現時点ではクイックソートと同程度であるとせざるを得ない。実験結果の検証を補強するために, 区分木ソートとヒープソートのcpu時間比を表に掲載した。表1から, クイックソートのヒープソートに対するcpu時間比が約1.9から5.9になることが読み取れる。D.R.Musserはクイックソートのヒープソートに対するcpu時間比が2から5になると指摘している [1]。表から読み取れるように, ヒープソートのクイックソートに対するcpu時間比は1.9から5.9の範囲に入っており, 本実験に用いたクイックソートやヒープソートの計算機実装方法, ならびに本実験結果はある程度の妥当性があると考えている。

実験により, k は約1.2とした時にcpu時間がもっとも少なかったが, その近辺の k に対しても性能の大きな劣化は見られなかった。

機能的な観点から見ると, バケット法や多重分割ソート, さらにはクイックソートなど多くのソートアルゴリズムがソート開始時点ですべてのキーが与えられていないのに対し, 区分木ソートではキーを逐次追加することができるという点と, stableであるという点を長所として持っている。この点はヒープソートも同様である。

5. 終わりに

作業メモリを $\Theta(n)$ だけ必要とはするものの, cpu時間が $O(n \log n)$ である2分木ソートを提案した。ソートの開始時期には2つの葉だけを持つ平衡2分木であったものが, キーが次々と追加されるにしたがって葉が2分割されて大規模な平衡2分探索木に生長するという特長を持つ。計算機実験により, 平均cpu時間がクイックソートと匹敵することを確認した。

[1] D.R.Musser, Introspective Sorting and Selection Algorithm, Software: Practice and Experiences,27(8),pp.987-pp.993, August,1997

[2] 川村畑行, 江口賢和, 小笠原基泰, 重村哲史, 多重分割ソートアルゴリズム, 情報処理学会, V o . 1 3 5 , N o . 1 2 , 1 9 9 4

[3] M.A.Bender, M.Farach-Colton, M.Mosteiro, Insertion Sort is $O(n \log n)$, http://xxx.lanl.gov/PS_cache/cs/pdf/0407/0407003.pdf

[4] S.Basse, アルゴリズム入門—設計と解析—, 星雲社, 1998

[付録] 比較に用いたクイックソートのアルゴリズム

```
Sub quicksort(key, a, z)
  If z - a <= 12 Then
    For i = a + 1 To z
      w=key(i)
      For j=i-1 to a step -1
        If key(j)>w then (key(j+1)=key(j)
                        else (key(j+1)=w: goto nexti )
      next j
      key(a)=w
nexti:  next i
  Else
    pivot = (key(a) + key((a + z) / 2) + key(z)) / 3
    left = a : right = z
    Do While left <= right
      Do While key(left) < pivot  (left = left + 1)
      Do While key(right) > pivot  (right = right - 1)
      If left <= right then (w = key(left): key(left) = key(right): key(right) = w:
                          left = left + 1 : right = right - 1)
    Call quicksort (key, a, right)
    Call quicksort (key, left, z)
  EndIf
```