

On Bounds of the Number of Comparisons to Select the
t-th Largest of n Elements

Kohei NOSHITA

(Department of Computer Science
University of Electro-Communications
Chofu, Tokyo 182)

Abstract

In this paper we shall show some results on the bounds of $V_t(n)$, where $V_t(n)$ denotes the number of binary comparisons necessary to select the t-th largest element of the given n elements.

Let $v(G,t)$ be the number of comparisons necessary to select the t-th largest element of G , where G is a partially ordered set of elements (poset, in short). The partial order is naturally defined as representing the orders among elements obtained from comparisons.

Let $\alpha(x,G)$ ($\beta(x,G)$) be the number of elements larger (smaller, respectively) than x in a poset G .

Theorem 1

Let G be a poset and G' be the poset obtained from G by deleting every element x such that either $\alpha(x,G) > t$ or $\beta(x,G) > n-t+1$. Let $S = \{x \mid \beta(x,G) > n-t+1\}$. Then $v(G,t) = v(G',t')$ for $t' = t - |S|$.

Theorem 2

$$V_t(n) \geq \min\{V_{t-1}(n-1), V_t(n-1)\} + 2$$
for $2 \leq t \leq \lceil n/2 \rceil$.

Theorem 3

$$V_5(10) \leq 16 .$$

Theorem 4 (Exact values of $V_t(n)$)

n	$V_1(n)$	$V_2(n)$	$V_3(n)$	$V_4(n)$	$V_5(n)$
6	5	7	8	8	7
7	6	8	10	10	10
8	7	9	11	12	12
9	8	11	12	(14)	14
10	9	12	(14)	(15)	(16)

The parenthesized values are the best upper bounds known until now. See the table on page 215 in the Knuth's book [2].

Theorem 5 (Improvement of the Hadian-Sobel algorithm)

There is a selection algorithm, which gives the following upper bound.

$$V_t(n) \leq n - t + (t-1) \lceil \log_2(n-t+2) \rceil - (t - \lfloor t/2 \rfloor - 1)$$

for n such that

$$2^{k+t-2} < n \leq 2^{k+t-2} + 2^{k - \lfloor t/2 \rfloor - 1}$$

where $k \geq 3$.

§1. はじめに

本稿では、選択アルゴリズム (selection algorithm) の基本的な性質を調べ、選択に要する比較回数¹の上下限に関するいくつかの結果を示す。

選択の問題とは、ソート²の問題の1種であり、任意に与えられた (大小の順序がつけられる) n 個の要素のうち、 t 番目に大きい要素を、 2 要素の大小比較 (binary comparison) の演算を用いて選択する (良い) アルゴリズムを見つけることである。本稿では特に、最悪の場合に要する比較の回数を最小にするアルゴリズムの考察に限定する。この問題は、既に Knuth の本 [2] にかなり詳しく最近の結果までふくめて紹介されているので、ここでも術語や従来³の結果は、特に断らない限り、同書を参照するものとする。

さて $V_t(n)$ で、与えられた n 個の要素のうち t 番目に大きい要素を求めるのに (最悪の場合に) 必要な比較回数 (の最小値) を表わすことにする。

本稿では、 $V_t(n)$ に関する基本的ないくつかの性質を調べ、その応用として、小さい n に対する $V_t(n)$ の正確な値を求め、また $V_5(10)$ に対する新しい上限値を与える。これは、[2] の 215 頁にて与えられている $V_t(n)$ の上限表のいくつかの項目の改良なしに正確な値を与えるものである。最後に $V_t(n)$

の上限を改良する1つのアルゴリズムを示す。

§2. 基本的な性質といくつかの $V_t(n)$ の値

入力として与えられる n 個の要素は, 1 から n までの番号と
ふつた n 個の数として考え, その番号で各要素を識別するも
のとする. 各要素の値は, 互へに相異なるものとしておく.

2 要素の比較によって判明した大小の順序は, 半順序とし
て考えることができるので, いくつかの要素でなるこの半順
序集合 (poset とよぶことにする) をグラフで図示することに
にする. これは, 普通の Hasse 図であるが, 間接的 (推察的)
な順序関係を示す枝は, 省略して描く. なお大きい要素は,
上に描くことにする. 特に入力 n 個の要素は, 孤立した n
個の点からなる poset である.

どんな選択アルゴリズムも, この poset を入力として, 比
較演算を用いて t 番目の要素と求める決定木 (decision tree) [2]
として表現できるが, 少し拡張解釈して, 1つの poset G を
入力として, G 中の 2 要素の比較を繰返して, t 番目の要素
と指定するものとする.

(Def 1)

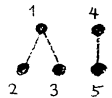
$v_A(G, t)$ は, poset G と順位 t に対する選択アルゴリズム
 ΔA が (最悪の場合に) 必要な比較回数とする.

(Def 2)

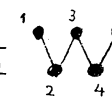
$$v(G, t) = \min_A \{v_A(G, t)\}.$$

定義より, $V_t(n) = v([\overset{\bullet}{1} \overset{\bullet}{2} \overset{\bullet}{3} \cdots \overset{\bullet}{n}], t).$

(Ex 1)

$$v([\overset{1}{\bullet} \overset{2}{\bullet} \overset{3}{\bullet} \overset{4}{\bullet} \overset{5}{\bullet}], 2) = 3.$$


(Ex 2)

$$v([\overset{1}{\bullet} \overset{2}{\bullet} \overset{3}{\bullet} \overset{4}{\bullet} \overset{5}{\bullet}], 3) = 3.$$


poset G に対して, $\#G$ でのその節点 (要素) の個数を表わすことにする. 便宜上, $t \leq 0$ または $t > \#G$ なる t に対して

$$v(G, t) = 0$$

としおく.

さて任意の poset G_1, G_2 がその要素の番号を除いて同じ形 (同型) なら, G_1 に対する選択アルゴリズムにおいて, 要素の番号に適当な置換を施せば, G_2 に対する一つの選択アルゴリズムに変換できる. 従って選択アルゴリズムや各定義は, 各 poset の節点に番号がふられていないと仮定しても本質的な差異はない.

(Def 3)

2つの poset G_1, G_2 に対して, 要素の間には 1対1 の対応 f があり, 任意の要素 $a, b \in G_1$ に対して,

$$a < b \text{ ならば } f(a) < f(b)$$

が G_2 で成立つ時,

$$G_1 \trianglelefteq G_2$$

とかく. (ここに $a \in G$ という記法で, a が G の要素であることと示す.)

(Def 4)

$$G_1 \equiv G_2 \quad \text{iff} \quad G_1 \trianglelefteq G_2 \quad \text{かつ} \quad G_2 \trianglelefteq G_1.$$

(Lemma 1)

任意の G_1, G_2, t に対して, $G_1 \trianglelefteq G_2$ ならば,

$$v(G_1, t) \geq v(G_2, t).$$

証明: $v(G_1, t)$ を実現するアルゴリズム A を simulate する (G_2, t) に対応するアルゴリズム B を考えればよい. ■

(Th 2)

任意の poset G と値 t に対して,

$$S_1 = \{ \nu \mid \#\{ \mu \mid \nu > \mu \} > n - t + 1 \}$$

$$S_2 = \{ \nu' \mid \#\{ \mu' \mid \nu' < \mu' \} > t \}$$

とする. いま G から S_1, S_2 に属するすべての要素を除いてできる poset を G' とすると,

$$v(G, t) = v(G', t').$$

ただし $t' = t - \#S_1$ とする.

証明: $v(G', t')$ を実現するアルゴリズム A とする.

A は、このまゝ (G, t) に対するアルゴリズム B にもなっているから、

$$v(G', t') \geq v_B(G, t) \geq v(G, t).$$

いま G の中で S_1 (S_2) に属す要素に (G の順序に矛盾することなく) 線形の順序をつけ、各要素の値は、 G' のどの要素よりも大きく (小さく) 与える "oracle" [2] をつくることができる。すると $v(G, t)$ を実現するアルゴリズム C は、 (G', t') に対するアルゴリズム D に変換できる。よって

$$v(G, t) \geq v_D(G', t') \geq v(G', t').$$

以上をあわせて $v(G, t) = v(G', t')$ が成立つ ■

なお別証明が [4] で与えられている。この定理は、特定の (G, t) に対するすべてのアルゴリズムを数えあげるとき、有用なものであり、ここでは特に "削除定理" とよぶことにする。

(Lemma 3°) [2]

$$V_t(n) \geq V_t(n-1) + 1$$

$$V_t(n) \geq V_{t-1}(n-1) + 1.$$

(Th 4)

$$V_t(n) \geq \min \{ V_{t-1}(n-1), V_t(n-1) \} + 2.$$

証明：つぎのいずれかが成立つことを証明する。

$$V_t(n) \geq V_{t-1}(n-1) + 2,$$

$$V_t(n) \geq V_t(n-1) + 2.$$

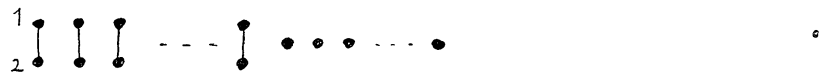
$V_t(n)$ を実現するアルゴリズム A を考える。 A は、次のような j 個の比較を始めるものと仮定してよい。

第 1 回目の比較: 要素 (1) と (2)

第 2 回目の比較: 要素 (3) と (4)

第 j 回目の比較: 要素 $(2j-1)$ と $(2j)$ 。

ただし $1 \leq j \leq \lfloor n/2 \rfloor$ 。これは、第 l 回目 ($1 \leq l \leq j$) の比較でできる 2 つの poset が同型であること、各比較に関わる要素が共通していないことによる (比較の結果できる 2 つの poset に依りて、 A が別の 2 要素を選ぶとしても、poset が同型であるから、 $V_t(n)$ を実現するアルゴリズム A' を作り、上の様に標準化できる)。さらに同様の考察により、 $(1) > (2)$ であると仮定してよい。ここまです、



のような poset になる。 A において、第 $(j+1)$ 回目の比較は、要素 (1) か (2) かのいずれかが一方の対象に選ばれ、もう一方は選ばれないと仮定することができる。いま A は、要素 (1) を選んだとする。この時要素 (1) に最大の値 (∞) を与える "oracle" をつくることができる。アルゴリズム \tilde{A} にお

て、要素 (1) に関する比較を全部除いてできるアルゴリズムを B とすると、 B は、 $(n-1)$ 個の要素のうち $(t-1)$ 番目の要素を選択するアルゴリズムになっている。アルゴリズム A は、要素 (1) について少なくとも 2 回比較しているから、

$$v_A(\overbrace{[\dots\dots\dots]}^n, t) - 2 \geq v_B(\overbrace{[\dots\dots\dots]}^{n-1}, (t-1)) \geq V_{t-1}(n-1)$$

$$\text{よって } V_t(n) \geq V_{t-1}(n-1) + 2.$$

もしアルゴリズム A が、上で要素 (2) を選んでいけば、最小の値を用いる同様の議論を用いると、

$$V_t(n) \geq V_t(n-1) + 2. \quad \blacksquare$$

(Ex 3)

$$V_3(7) = 10.$$

Floydによつて、この事実はすでに報告されている [6]。ここでは、(lemma 1) と (削除定理) の応用の例として、(7, 3) に対するオベアのアルゴリズムを表現する minimax 法に基づく tree をしらみつぶして証明した。これは、かなり大きなサイズの "tree search" であるが、強力な "定理" の援用で "手" で計算できる範囲におとすことができた 1 例である。"計算機" でこの tree search を実行させるのは、容易ではない。なぜなら、単純な方法では、グラフの同型性判定の (高効率な) アルゴリズムが必要になるからである。

(Cor 5)

$$V_3(8) = 11 .$$

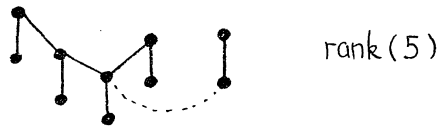
$V_5(10)$ に関しては, [2] によると $V_5(10) \leq 17$ となる。
 次はこの上限を改良する。

(Th. 6)

$$V_5(10) \leq 16 .$$

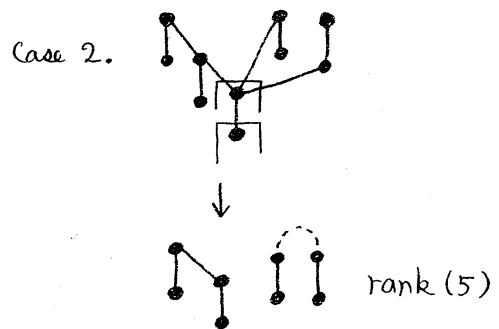
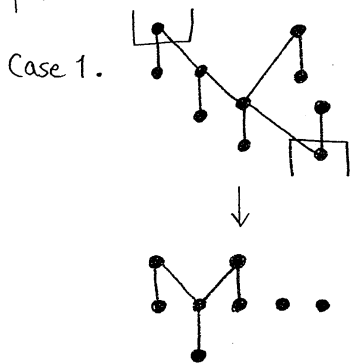
証明: 16 回の比較で選択する一つのアルゴリズムを示す。

Step 8 の直後:

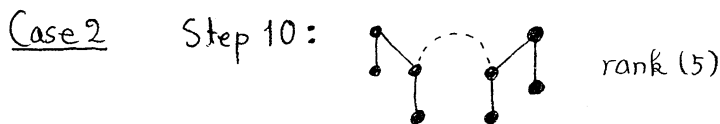


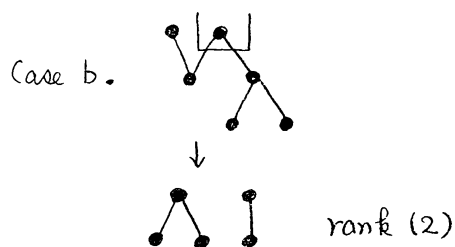
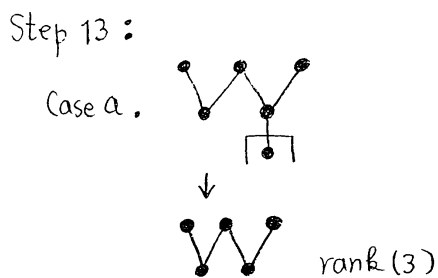
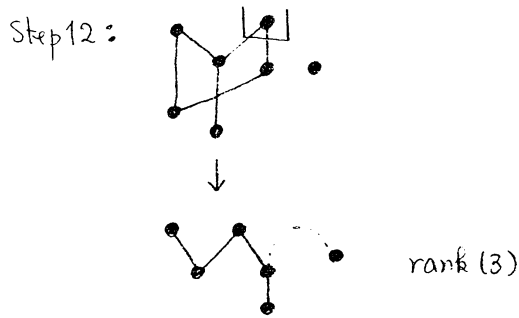
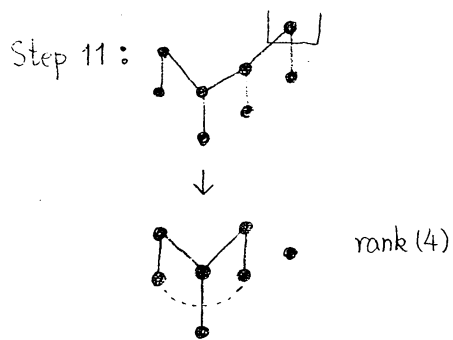
(8 回の比較でここまで (対称的に) くることが出来るのは, 明らかである.)

Step 9:



Case 1. については, あと 7 回の比較で十分である ([3] に
 おける証明を参照)。





OK (cf. (Ex 2))

OK (cf. (Ex 1))

いずれの場合も 16 ステップで選択できる。 ■

(Th 7)

小さい n に対する $V_t(n)$ の値は、次の通りである。

	$V_1(n)$	$V_2(n)$	$V_3(n)$	$V_4(n)$	$V_5(n)$	$V_6(n)$
6	5	7	8	8	7	5
7	6	8	10*	10*	10*	8
8	7	9	11*	12	12	11*
9	8	11	12	(14)	14	14
10	9	12	(14)	(15)	(16)	(16)

ここで () で括弧しているものは、いままで知られていない最善の上限である。太字が本稿で求めた結果である。印*の

ついでに t のについては, (Ex 3) の注釈参照. その他の t は, 既に正確な値として知られている [2]. なお $V_3(9) = 12$ は, [5] による.

$V_5(9) \leq 14$ は, [3] による. これが正確な値であることは, $V_4(8) = 12$ とともに, (Th 4) による.

$V_4(9) = 13$ または 14 である (14らしい). $V_4(9) = 14$ が確定すると, $V_5(10) = 16$ も確定する (Th 4).

さて, ここで述べた一連の結果は, 一般に, 中間値 (median) を求める問題について,

$$V_{\lceil n/2 \rceil}(n) = 2n - 4$$

と予想がたてられそうであるが, 現在知られている上限と下限との間隙は, まだ相当大きい.

§3. Hadian-Sobel アルゴリズムの改良

任意の n, t に対して, n に線型的な $V_t(n)$ の上限を与えるアルゴリズムがよく知られているが [1], 順位 t が小さい時には, Hadian-Sobel のアルゴリズムがはるかに能率的であることも知られている [2]. 本節では, この Hadian-Sobel のアルゴリズムを改良する.

(Th 8)

任意の n, t ($3 \leq t \leq \lceil n/2 \rceil$) に対して,

$$2^k + (t-2) < n \leq 2^k + (t-2) + 2^{k - \lfloor t/2 \rfloor - 1}$$

ならば ($k \geq 3$),

$$V_t(n) \leq (n-t) + (t-1) \lceil \log_2(m-t+2) \rceil - (t - \lfloor t/2 \rfloor - 1)$$

である。

上の範囲の n に対し、Hadian-Sobel のアルゴリズムが与える上限より、

$$(t - \lfloor t/2 \rfloor - 1) \text{ 回}$$

改良されることになる。他の n については、そのままだである。

証明: 上の上限式と与える 1 つのアルゴリズムと与える。

$r = \lfloor t/2 \rfloor + 1$ とおく。また $n = 2^k + (t-2) + m$ (ただし $1 \leq m \leq 2^{k-r}$) とする。

(1) n 個中 $(t-2)$ 個の要素を (Hadian-Sobel のアルゴリズムと同様に) 最後のステップのためにとっておく。

(2) 残りの $(n-t+2)$ 個の要素を (2^r+1) 個の集合 $S_1, S_2, \dots, S_{2^r}, R$ に分割する。ここで S_i ($i=1, \dots, 2^r$) は、 2^{k-r} 個、 R は、 m 個の要素をもつものとする。

(S_i の要素の個数は、2 の中であることに注意)

[Step 1] 各 S_i ($i=1, \dots, 2^r$) に対し、標準的な (バランスしてやる) knockout tournament を実施して、最大の要素 σ_i と選択する。

いま $M_1 = \{\sigma_1, \sigma_2, \dots, \sigma_{2^r-1}\}$, $M_2 = \{\sigma_{2^r+1}, \dots, \sigma_{2^r}\}$ とおく。

[Step 2] 各 M_i ($i=1, 2$) に対して, 標準的な knockout tournament を実施して, 最大の要素と選択する.

ここで, 次のような M_i の要素と σ_{j_i} とする. [Step 2] の tournament において, σ_{j_i} は, 第 1 回戦で負けた (小さい) 要素であり, その比較で勝った要素は, 第 2 回戦で負けた要素であり, ..., その比較で勝った要素は, 最終戦で負けた要素である, ようなものである. たとえば図 1 では, \circ で囲んである要素が各々 σ_{j_i} である. ここで $1 \leq j_1 \leq 2^{r-1}$, $(2^{r-1} + 1) \leq j_2 \leq 2^r$ である.

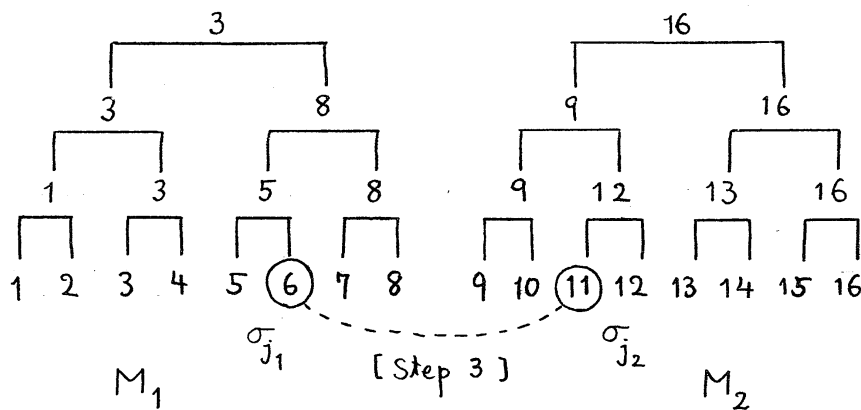


図 1. Step 2 における knockout tournament ($r=4$)

[Step 3] σ_{j_1} と σ_{j_2} と比較する.

対称性より, [Step 3] で σ_{j_2} が負けた要素と仮定して差支えない. 各 σ_{j_i} ($i=1, 2$) は, 自分より大きい要素が, 少なくとも $\lfloor t/2 \rfloor$ 個あるから, σ_{j_2} は, 決して t 番目に大

きい要素になれない。したがって S_{j_2} の要素をすべて除くことができる。

[Step 4] S_{j_2} のかわりに R を考える。まず R において、標準的な knockout tournament を実施し、最大の要素と選択する。[Step 2] で実施した M_2 に対する tournament で残っている要素の集合において、いま求めた R の最大の要素に関する置換え選択 (replacement selection) を実施する。これで選択された最大の要素と ([Step 2] で既に得られている) M_1 の最大の要素と比較する。

この比較で大きい方の要素は、この段階まで考慮した $(2^k + m)$ 個の要素の最大要素であるから、除くことができる。

[Step 5] (これ以後は Hadian-Sobel アルゴリズムと同様) 最初にとつておいた $(t-2)$ 個の要素を一つずつ取り上げ、置換え選択を実施して、 $(t-2)$ 個の要素を除く。最後に、残った要素の中から最大の要素を選べば、それが望みの t 番目に大きい要素である。

以上でアルゴリズムの記述は、おしまいである。このアルゴリズムが正しいことは、容易に検証できる。

次にこのアルゴリズムで要する比較回数を計算する。

$$[\text{Step 1}] \quad 2^r \times (2^{k-r} - 1) \quad \text{回}$$

$$[\text{Step 2}] \quad 2^r - 2 \quad \text{回}$$

$$[\text{Step 3}] \quad 1 \quad \text{回}$$

$$[\text{Step 4}] \quad (\text{高々}) \quad (m-1) + (r-1) + 1 \quad \text{回}$$

$$[\text{Step 5}] \quad (\text{高々}) \quad (t-2) \times k + (k-1) \quad \text{回}$$

併せて定理で述べている $V_t(n)$ の上限式を得る。 ■

特に $t=3$ とすると,

$$2^k + 1 < n \leq 2^k + 2^{k-2} + 1$$

とみたす n に対して,

$$V_3(n) \leq n - 4 + 2^{\lceil \log_2(n-1) \rceil}$$

となり, [5] の下限の結果とあわせて, 上のよ様な n に対して, $V_3(n)$ の上下限の差は, 1回となる (次頁 (Note) 参照).

§4. あとがき

本稿では, 選択アルゴリズムのいくつかの基本的な性質を調べ, 小さい n に対する $V_t(n)$ の正確な値を求めた. また

$V_5(10) \leq 16$ と示した. こからの結果は, Knuth の本 [2] の $V_t(n)$ の表の改良と与えている. 最後に, Hadian-Sobel のアルゴリズムの1つの改良を示した.

(Note)

最新の Knuth からの私信 [6] には, 次の事実が伝えられてゐる。

- (1) D. Kirkpatrick が, すべての $n \geq 50$ に対して, $V_3(n)$ の正確な値と決定したこと。
- (2) 同じく, §4. で述べた Hadion-Sobel のアルゴリズムの改良も独立に示したこと, 比較回数の式は, 本稿のものと同じであるが, さらに n の範囲に対する制限を弱めたこと。

§5. 引用文献

- [1] Blum, M. et al, "Time Bound for Selection," JCSS, Vol. 7 (1973), 448-461.
- [2] Knuth, D.E., The Art of Computer Programming, Vol. 3 (Sorting and Searching), Addison-Wesley (1973).
- [3] Noshita, K., "Median Selection of 9 Elements in 14 Comparisons," Info. Proc. Letters, Vol. 3, No. 1 (1974), 8-12.
- [4] 野崎昭弘, 私信 (1975年3月).
- [5] Pratt, V. & F.F. Yao, "On Lower Bounds for Computing The i -th Largest Element," IEEE 14th SWAT (1973), 70-81.
- [6] Knuth, D.E., 私信 (1975年5月). ■