Implementations of Kruskal algorithm

Giovanni Righini



Median-of-medians

Median selection problem: find the median in a given set (unsorted list *L*) of *n* values.

Median-of-medians algorithm.

Selection

- Split *L* in subsets of 5 elements each: $S_i \ \forall i = 1, ..., n/5$.
- Find the median x_i for each S_i.
- Find the median M of the x_i values.
- Partition L into three lists:
 - L₁ with values smaller than M;
 - L₂ with values equal to M;
 - L₃ with values larger than M.
- Discard L₁ or L₃ depending on the sought element k.
- Recursively call the procedure on the remaning part of *L*.



```
// Procedure Select(L, u, v, k).
if (v - u + 1 < 10) then
  Sort(L, u, v)
   return L[u-1+k]
else
   s \leftarrow \lceil (v - u + 1)/5 \rceil
  for i = 1, \ldots, s do
     S(i) \leftarrow \{u + 5(i - 1), \dots, \min\{u + 5i - 1, v\}\}
     x(i) \leftarrow \text{Median} - \text{of} - 5(S(i))
   M \leftarrow \text{Median}(\{x(1), \dots, x(s)\})
   Partition(L, M, L_1, L_2, L_3)
   if (k < |L_1|) then
     Select(L_1, u, u + |L_1| - 1, k)
  else
     if (k > |L_1| + |L_2|) then
        Select(L_3, v - |L_3| + 1, v, k - |L_1| - |L_2|)
     else
        return M
```



Complexity

Median of 5: at most 6 pair-wise comparisons.

Median of a set *S* of *s* elements: call Select(S, 1, s, $\lceil s/2 \rceil$).

Size s of S is n/5.

Selection

By definition, at least half of the medians x_i (i.e. at least n/10) are $\leq M \, (\geq M)$.

For each of them, 3 values in S_i are $\leq M$ ($\geq M$).

Then, $|L_1| \ge 3n/10$ ($|L_3| \ge 3n/10$).

Hence $|L_3| \le 7n/10$ ($|L_1| \le 7n/10$).



Complexity

T(n): n. of comparisons to find the median of n elements.

Finding all medians x_i takes at most 6n/5 comparisons.

Finding *M* takes T(n/5).

Partitioning *L* takes n-1 comparisons.

Finding the median of the remaining list takes at most T(7n/10).

Therefore

Selection

$$T(n) \leq \frac{6n}{5} + T\left(\frac{n}{5}\right) + n + T\left(\frac{7n}{10}\right).$$

Hence

$$T(n) \leq \frac{11n}{5} \sum_{k=0}^{\infty} \left(\frac{1}{5} + \frac{7}{10}\right)^k.$$

Since $\frac{1}{5} + \frac{7}{10} = \frac{9}{10} < 1$, the geometric series converges to a constant.

Therefore T(n) is O(n) (it can be reduced to 2.95 n).



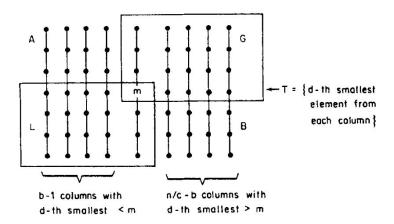
Selection problem

Selection problem: find the k^{th} smallest value in a given set (unsorted list N) of n values.

Select algorithm (Blum et al., 1973).

- Split N in n/c subsets of c elements each and sort each of them.
- Find the dth smallest element (the representative) in each subset.
- Find the b^{th} smallest element, M, in the set of representatives.
- Partition the n/c subsets into three parts:
 - Left: b − 1 subsets with representatives < M;
 - Center: the subset whose representative is M;
 - Right: n/c b subsets with representatives $\geq M$.
- L: entries in the left and center columns on rows from 1 to d.
- G: entries in the right and center columns on rows from d to c.
- Compare each element with *M*.
- Discard either L or G and repeat, until k − 1 smallest entries or n-k largest entries have been discarded.

Selection problem





Selection

Complexity

h(c): n. of comparisons needed to sort c elements.

Ford and Johnson (1959): $h(c) = \sum_{i=1}^{c} \lceil \log_2(\frac{3}{4}j) \rceil$.

T(n): n. of comparisons needed to find the k^{th} smallest element in a set of n.

- Sorting $\frac{n}{c}$ subsets of c elements each: $\frac{n}{c}h(c)$.
- Finding M among the representatives: $T(\frac{n}{2})$.
- Comparing each element with M: n.
- Solving the remaining problem: $T(n \min\{|L|, |G|\})$.

$$T(n) = \frac{n}{c}h(c) + T(n/c) + n + T(n - \min\{|L|, |G|\}).$$



$$|L| = bd$$
 $|G| = (\frac{n}{c} - b + 1)(c - d + 1)$

By selecting b = n/42, c = 21, d = 11:

$$|L| = \frac{11}{42}n$$
 $|G| = \left(\frac{n}{21} - \frac{n}{42} + 1\right)11 = \frac{11}{42}n + 11.$
 $h(21) = 66.$

Therefore

Selection

$$T(n) \leq \frac{66}{21}n + T\left(\frac{n}{21}\right) + n + T\left(\frac{31}{42}n\right).$$

Base of the induction: for small enough n, h(n) is linear (e.g. for $n < 10^5$, h(c) < 19n).

$$T(n) \leq \frac{29}{7}n\sum_{k=0}^{\infty} \left(\frac{1}{21} + \frac{31}{42}\right)^k = \frac{29}{7}n\frac{1}{1 - \frac{11}{14}} = \frac{29}{7}n\frac{14}{3} = \frac{58}{3}n.$$

Hence T(n) is O(n).



Partial sorting problem

Partial sorting problem: find and sort the k^{th} smallest values in a given set N of m values.

SelectionSort algorithm.

- For i = 1, ..., k.
 - scan the list from N[i] and find the smallest value;
 - swap the smallest element with N[i].

Time complexity: O(km). Space complexity: O(m).

Using a (big) binary heap (C++ STL's partial_sort function):

- Build a binary heap with the m elements.
- For k times, extract the root and rearrange the heap.

Complexity: $O(m + k \log m)$. Space complexity: O(m).



Partial sorting problem

Using a (small) binary heap:

- Build a binary max-heap with the first k elements of N.
- For i = k + 1, ..., m times, if N[i] is smaller than the current root of the heap, then replace the root and update the heap.
- Sort the elements of the heap.

Complexity: $O(k + m \log k)$. Space complexity: O(k). Valid alternative if k << m or N is processed on-line.

QuickSelSort algorithm, combining Select and QuickSort.

- find the kth smallest element in N with Select (the pivot);
- scan the list and retain the elements smaller than the pivot (if needed);
- sort the selected elements.

Complexity: $O(m + k \log k)$. Space complexity: O(m).



PartialQuickSort (Martinez, 2004)

Input for each recursive call of PartialQuickSort(N, i, j, k):

- an array N;
- an interval $[i, \ldots, j]$, with $i \leq j$;
- a number k of smallest elements of N to sort, with i < k.

Base of the recursion: i = j.

Recursive step:

- select a pivot in position p ∈ [i, ..., j];
- partition $[i, \ldots, j]$ as in QuickSort; let p' be the final position of the pivot:
- recursively call PartialQuickSort(N, i, p' − 1, k);
- if k > p', then recursively call PartialQuickSort(N, p' + 1, j, k).



PartialQuickSort

```
// Procedure PartialQuickSort(N, i, j, k)

if (i < j) then
p \leftarrow \text{Pivot}(N, i, j)
p' \leftarrow \text{Partition}(N, i, j, p)
PartialQuickSort(N, i, p' - 1, k)

if (p' < k) then
PartialQuickSort(N, p' + 1, j, k)
```



Incremental sorting problem (Paredes and Navarro, 2006)

Incremental sorting problem: on-line version of the Partial sorting problem (*k* is not given).

Given set N of m numbers, output the elements of N from the smallest one to the largest one, so that the process can be stopped after k elements have been output, for any k (not given in input).

Iterate Select.

Repeatedly select the smallest unselected element in N.

Time complexity: O(km). Space complexity: O(m).

Using a (big) binary heap:

- Build a binary heap with the m elements.
- Repeatedly, extract the root and rearrange the heap.

Complexity: $O(m + k \log m)$. Space complexity: O(m).

These naive methods are dominated by PartialQuickSort.



Incremental Quick Select

Algorithm IncrementalQuickSelect (Paredes and Navarro, 2006).

At each iteration, the algorithm finds the smallest element among those not yet selected.

The algorithm exploits the same recursion of QuickSort, but it can stop for each value of k.

As in iterative implementations of QuickSort, it keeps a stack of pivots already set at their correct positions in previous iterations.



Incremental Quick Select

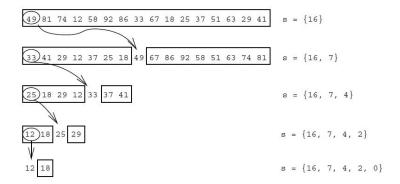


Figure: Positions are numbered from 0 to 15.



Incremental Quick Select

Incremental sorting

Algorithm 1 Procedure IQS(N)

```
1: S \leftarrow \{s\}
2: k \leftarrow 0
 3: repeat
       if Top(S) = k then
         Pop(S)
 5:
 6: return N[k]
 7: k \leftarrow k + 1
 8: else
         p \leftarrow \mathsf{Pivot}(N, k, \mathsf{Top}(S) - 1)
 9:
          p' \leftarrow \text{Partition}(N, k, \text{Top}(S) - 1, p)
10:
          Push(S, p')
11:
12: until k = s
```

Worst-case complexity: $O(m^2)$. Average-case complexity: $O(m + k \log k)$.



Comparing MST algorithms

Computational complexity:

- Fibonacci heaps (Fredman and Tarjan, 1987, Gabow et al. 1986) for Prim algorithm: $O(m + n \log n)$.
- Union-Find for Kruskal algorithm: O(m log n) for sorting and O(m + n log n) for the MST.

A randomized algorithm by Karger (1993) runs in linear expected time.

In practice:

- Kruskal: effective when the edge weights are drawn from a small range of integers or when the graph is sparse;
- Prim: effective in all the other cases, especially when the graph is dense.

Experimental comparisons (Moret and Shapiro, 1991, Bazlamaçcı and Hindi, 2001):

- Prim algorithm implemented with a binary heap $(O(m \log n))$ is faster than Kruskal.
- Both algorithms are faster than the more recent ones.



Implementing Kruskal algorithm

Naive implementation: first the edge list is sorted; then it is scanned to compute a MST.

"On demand sorting": keep the edges partially sorted in a binary heap, constructed by QuickSort when needed, i.e. when all previously sorted edges have been examined.

QuickSortKruskal (QSK): no binary heap is constructed. Kruskal algorithm is executed from within QuickSort: the two steps of sorting and selecting the edges of the MST proceed in parallel.



Quick Sort Kruskal

QSK: the recursive procedure is run on the leftmost part first and on the rightmost part later, for each interval to be sorted.

Base of the recursion: interval with a single element.

Every time QuickSort completes the sorting of the leftmost part of an interval, the pivot element is considered for insertion in the MST as in Kruskal algorithm, before sorting the rightmost part of the interval.

- edges are sorted only when needed;
- no additional data-structure and initialization.

QSK is especially effective when the input graph is given as an edge list.



Quick Sort Kruskal

```
// QuickSortKruskal. IN: \mathcal{E}. OUT: \mathcal{T}. InitUnionFind \mathcal{T} \leftarrow \emptyset count \leftarrow 0 QuickSortK(1, m)
```

 \mathcal{E} : input list of edges. \mathcal{T} : list of selected edges.

count: cardinality of T.



// Procedure QuickSortK(p, q)

```
e^+ \leftarrow a
if p < q then
   e^- \leftarrow p
   while e^- < e^+ do
      while \mathcal{E}[e^+].cost > \mathcal{E}[p].cost do
         e^{+} \leftarrow e^{+} - 1
      while (e^- \le e^+) \wedge (\mathcal{E}[e^-].cost \le \mathcal{E}[p].cost) do
         e^- \leftarrow e^- + 1
      if e^- < e^+ then
         Swap(e^-, e^+); e^- \leftarrow e^- + 1; e^+ \leftarrow e^+ - 1
   Swap(p, e^+)
if e^+ > p then
   QuickSortK(p, e^+ - 1)
if count < n-1 then
   TestEdge(e<sup>+</sup>)
if (count < n - 1) \land (e^+ < q) then
   QuickSortK(e^+ + 1, q)
```



Quick Sort Kruskal

The pivot element in position e^+ is tested by $TestEdge(e^+)$

- after sorting the edges with cost smaller than the pivot (QuickSortK(p, e⁺ - 1))
- before sorting the edges with cost larger than the pivot (QuickSortK(e⁺ + 1, q)).

```
// Procedure TestEdge(e)
i \leftarrow \mathcal{E}[e].i
j \leftarrow \mathcal{E}[e].j
if head[i] \neq head[j] then
count \leftarrow count + 1
T \leftarrow T \cup \{[i,j]\}
if card[head[i]] \geq card[head[j]] then
Append(head[j], head[i])
else
Append(head[i], head[j])
```



The Star Quick Sort Kruskal algorithm

Input graph: set of stars.

Naive solution 1: produce a sorted list of edges by merging the subsets and then sorting the resulting list or heap.

Naive solution 2: produce a sorted list of edges by separately sorting each star and then merging them into a unique sorted list or heap.

In StarQuickSortKruskal (SQSK) stars are not merged.

A sorted list is produced on demand by QuickSort from the edges of each star; the first not-yet-examined edge in each list is the candidate edge for its vertex.

All candidate edges are partially sorted in a binary heap.

Iteratively:

- the best candidate e is selected;
- QuickSort is re-activated to find the next candidate for the endpoints of e;
- the heap is rearranged.



The Star Quick Sort Kruskal algorithm

Main idea: save a potentially large fraction of the computing time that would be spent to (partially) sort each vertex star in a binary heap or a in sorted list, since only a very small fraction of the edges in each star is likely to be considered by Kruskal algorithm.

The candidate edge in each star is the last edge of the sorted part of the star; all the edges with a cost larger than the candidate remain unsorted.

Pro: the sorting step does not work on a (typically large) set of edges, but on several (much smaller) subsets.

Drawback: the selection of each edge requires some non-trivial steps.

Since sorting requires super-linear computing time, it is intuitively convenient to (partially) sort n subsets of cardinality m/n rather than a unique set of cardinality m, especially when the input is given as a set of vertex stars.



SQSK: iterative Quick Sort

Since in SQSK QuickSort must be executed step-by-step only on-demand, an iterative implementation is needed.

A step of QuickSort is executed by QSstep(v) on Star(v) every time v is the endpoint of the selected edge.

The effect of QSstep is to define at least one more sorted element in Star(v).

Stack(v): a stack associated with the star of each vertex v.

The stack contains the intervals in which the star has been partitioned by QuickSort.

Every time the current interval [p, ..., q] is divided into a left interval [p, ..., k-1] and a right interval [k+1, ..., q] by a pivot element k,

- the right interval [k+1,...,q] is put into the stack,
- the pivot element k is put in its final position between the two intervals,
- the left interval [p, ..., k-1] becomes the new current range and become the new current range.

SQSK: iterative QuickSort

h(v): position of last sorted element in Star(v).

If h(v) + 1 coincides with the leftmost element of the interval on top of Stack(v), then the new candidate edge is already available: increase h(v) and stop.

Otherwise:

- extract [p, q] to be sorted from Stack(v);
- select a pivot element k in [p, q];
- temporarily put the pivot in the first position;
- separate all elements smaller than or equal to the pivot (left) from the others (right);
- put the pivot back between the two sub-intervals;
- store the right interval [k + 1, q] in Stack(v), unless its cardinality is less than 2;
- repeat on the left interval [p, k-1], until its cardinality is less than 2.

After this step at least one more edge has been sorted and h(v) can be increased by 1, indicating the next candidate in Star(v).

```
// Procedure QSstep(v)
```

if
$$(Stack(v) \neq nil) \land (h(v) + 1 = Top(v).p)$$
 then $[p,q] \leftarrow Pop(v)$ while $p < q$ do $k \leftarrow Pivot(v,p,q)$ $Swap(v,k,p); e^+ \leftarrow q; e^- \leftarrow p+1$ while $e^- \leq e^+$ do while $Star(v)[e^+].cost > Star(v)[p].cost$ do $e^+ \leftarrow e^+ - 1$ while $(e^- \leq e^+) \land (Star(v)[e^-].cost \leq Star(v)[p].cost)$ do $e^- \leftarrow e^- + 1$ if $e^- < e^+$ then $Swap(v,e^-,e^+); e^- \leftarrow e^- + 1; e^+ \leftarrow e^+ - 1$ if $(e^+ \neq p)$ then $Swap(v,p,e^+)$ if $e^+ + 1 < q$ then $Push(v,e^+ + 1,q)$ $q \leftarrow e^+ - 1$ $h(v) \leftarrow h(v) + 1$



SQSK: initialization

```
// Procedure InitSQSK for i = 1, ..., n do Stack(i) \leftarrow \emptyset Push(i, [1, |Star(i)|]) h(i) \leftarrow 0 QSstep(i) BuildHeap(Heap) T \leftarrow \emptyset count \leftarrow 0 InitUnionFind
```



```
// Procedure SQSK, IN: Star. OUT: T.
InitSQSK
while count < n-1 do
  i \leftarrow Heap[1]
  i \leftarrow Star(i)[h(i)].vertex
  w \leftarrow Star(i)[h(i)].cost
  QSstep(i)
  (QSstep(i))
  IncreaseKev(i)
  (IncreaseKev(i))
  if head[i] \neq head[j] then
     T \leftarrow T \cup \{[i, i]\}
     count \leftarrow count + 1
     UpdateUnionFind(i, i)
```

Parenthesized instructions are needed if there are duplicates (each object to two stars).

The pivot selection

The selection of the pivot is crucial role for the performance of all algorithms (like QSstep) that need to partition intervals recursively.

QuickSort.

Purpose: completely sort a set of values.

Pivot selection: try to achieve a balanced partition.

SQSK.

Purpose: quickly compute the next candidate edge when required.

Pivot selection: go for an unbalanced partition.

Most part of each star is likely to be useless: the initial left interval should better be much smaller than the initial right interval.



The pivot selection

Heuristic rule: select r elements at random with uniform probability distribution in $[p, \ldots, q]$ and take the min cost one as the pivot element.

The larger r, the more unbalanced the resulting partition is likely to be.

Rule of thumb: set *r* to min $\{(q-p+1)/100, \overline{r}\}$.

Set \overline{r} to a larger value in InitSQSK and to smaller values otherwise.

There is room for heuristics, especially self-adaptive heuristics.



Experimental results

QSK: an iteration of Kruskal algorithm is done sometimes during the execution of QuickSort.

SQSK: an iteration of QuickSort is done sometimes during the execution of Kruskal algorithm.

QSK is designed to work on a list of edges.

SQSK is designed to work on a set of vertex stars.

Both QSK and SQSK can be faster than Prim algorithm.

Besides size and density, the computing time is also affected by the clustering degree of the input graph:

Non-clustered: SQSK wins.

Clustered: Prim wins.

