Advanced Topics in Sorting

complexity
system sorts
duplicate keys
comparators

▶ complexity

system sorts
 duplicate keys
 comparators

Computational complexity. Framework to study efficiency of algorithms for solving a particular problem X. Machine model. Focus on fundamental operations.

Upper bound. Cost guarantee provided by some algorithm for X. Lower bound. Proven limit on cost guarantee of any algorithm for X. Optimal algorithm. Algorithm with best cost guarantee for X.



Example: sorting.

Complexity of sorting

- Machine model = # comparisons <---- access information only through compares
- Upper bound = N lg N from mergesort.
- Lower bound ?



Comparison-based lower bound for sorting

Theorem. Any comparison based sorting algorithm must use more than N lg N - 1.44 N comparisons in the worst-case.

Pf.

- Assume input consists of N distinct values a1 through aN.
- Worst case dictated by tree height h.
- N! different orderings.
- (At least) one leaf corresponds to each ordering.
- Binary tree with N! leaves cannot have height less than Ig (N!)

h ≥ lg N!



≥ lg (N / e) ^N ← Stirling's formula

= N lg N - N lg e

≥ N lg N - 1.44 N

Complexity of sorting

Upper bound. Cost guarantee provided by some algorithm for X. Lower bound. Proven limit on cost guarantee of any algorithm for X. Optimal algorithm. Algorithm with best cost guarantee for X.

Example: sorting.

- Machine model = # comparisons
- Upper bound = N lg N (mergesort)
- Lower bound = N lg N 1.44 N

Mergesort is optimal (to within a small additive factor)

lower bound = upper bound

First goal of algorithm design: optimal algorithms

Complexity of sorting in context

Mergesort is optimal (to within a small additive factor)

Other operations?

- statement is only about number of compares
- quicksort is faster than mergesort (lower use of other operations)

Space?

- mergesort is not optimal with respect to space usage
- · insertion sort, selection sort, shellsort, quicksort are space-optimal
- is there an algorithm that is both time- and space-optimal?

stay tuned for heapsort

Nonoptimal algorithms may be better in practice

- statement is only about guaranteed worst-case performance
- quicksort's probabilistic guarantee is just as good in practice

Lessons

don't try to design an algorithm that uses half as many compares as mergesort

- use theory as a guide
- know your algorithms
- use guicksort when time and space are critical

Complexity of sorting in context (continued)

Lower bound may not hold if the algorithm has information about

- the key values
- their initial arrangement

Partially ordered arrays. Depending on the initial order of the input, we may not need N lg N compares.

 insertion sort requires O(N) compares on an already sorted array

Duplicate keys. Depending on the input distribution of duplicates, we may not need N lg N compares.

stay tuned for 3-way quicksort

Digital properties of keys. We can use digit/character comparisons instead of key comparisons for numbers and strings.



Another complexity example: Selection Selection: guick-select algorithm Find the kth largest element. Partition array so that: if k is here if k is here • Min: k = 1. • element a [m] is in place set r to m-1 set | to m+1 • Max: k = N. • no larger element to the left of m • Median: k = N/2. • no smaller element to the right of m Repeat in one subarray, depending on m. 1 r Applications. Order statistics. Finished when m = k \leftarrow a[k] is in place, no larger element to the left, no smaller element to the right Find the "top k" haddindulu tableten deductent Use theory as a guide public static void select(Comparable[] a, int k) հանվեսինի հմինքերը հայինքի 🖩 • easy O(N log N) upper bound: sort, return a[k] StdRandom.shuffle(a); վերիկինին, երկելին • easy O(N) upper bound for some k: min, max int 1 = 0;• easy $\Omega(N)$ lower bound: must examine every element int r = a.length - 1;while (r > 1)int i = partition(a, l, r); Which is true? and the state of the if (m > k) r = m - 1;• Ω(N log N) lower bound? [is selection as hard as sorting?] else if (m < k) l = m + 1;else return; • O(N) upper bound? [linear algorithm for all k] } 10

11

Quick-select analysis

Theorem. Quick-select takes linear time on average.

Pf.

- Intuitively, each partitioning step roughly splits array in half.
- N + N/2 + N/4 + ... + 1 \approx 2N comparisons.
- Formal analysis similar to quicksort analysis:

$C_{N} = 2 N + k \ln (N / k) + (N - k) \ln (N / (N - k))$

Ex: (2 + 2 In 2) N comparisons to find the median

Note. Might use ${\sim}N^2/2$ comparisons, but as with quicksort, the random shuffle provides a probabilistic guarantee.

Theorem. [Blum, Floyd, Pratt, Rivest, Tarjan, 1973] There exists a selection algorithm that take linear time in the worst case. Note. Algorithm is far too complicated to be useful in practice.

Use theory as a guide

- still worthwhile to seek practical linear-time (worst-case) algorithm
- until one is discovered, use quick-select if you don't need a full sort



Sorting Challenge 1 Problem: sort a file of huge records with tiny keys. Ex: reorganizing your MP3 files. Which sorting method to use? 1. mergesort 2. insertion sort 3. selection sort 101 Brown 1 A 243-456-9091 Fox file 📥 Quilici 1 C 343-987-5642 32 McCosh Chen 2 A 884-232-5341 11 Dickinson Furia 3 A 766-093-9873 22 Brown Kanaga 3 в 898-122-9643 343 Forbes record Andrews 3 A 874-088-1212 121 Whitman 232-343-5555 115 Holder Rohde 3 Battle 4 C 991-878-4944 308 Blair 4 A 664-480-0023 097 Little Aaron 4 B 665-303-0266 113 Walker Gazsi

Sorting Challenge 1

Problem: sort a file of huge records with tiny keys. Ex: reorganizing your MP3 files.

Which sorting method to use?

- 1. mergesort probably no, selection sort simpler and faster
- 2. insertion sort no, too many exchanges
- 3. selection sort YES, linear time under reasonable assumptions

14

16

Ex: 5,000 records, each 2 million bytes with 100-byte keys.

- Cost of comparisons: 100 × 5000² / 2 = 1.25 billion
- Cost of exchanges: 2,000,000 × 5,000 = 10 trillion
- Mergesort might be a factor of log (5000) slower.

Sorting Challenge 2

Problem: sort a huge randomly-ordered file of small records. Ex: process transaction records for a phone company.

Which sorting method to use?

- 1. quicksort
- 2. insertion sort
- 3. selection sort

file 🔺	Fox	1	A	243-456-9091	101 Brown
	Quilici	1	с	343-987-5642	32 McCosh
	Chen	2	A	884-232-5341	11 Dickinson
	Furia	3	A	766-093-9873	22 Brown
	Kanaga	3	в	898-122-9643	343 Forbes
record 📥	Andrews	3	A	874-088-1212	121 Whitman
	Rohde	3	A	232-343-5555	115 Holder
	Battle	4	с	991-878-4944	308 Blair
kev 📥	Aaron	4	A	664-480-0023	097 Little
	Gazsi	4	в	665-303-0266	113 Walker

Sorting Challenge 2

13

15

Problem: sort a huge randomly-ordered file of small records. Ex: process transaction records for a phone company.

Which sorting method to use?

- 1. quicksort YES, it's designed for this problem
- 2. insertion sort no, quadratic time for randomly-ordered files
- 3. selection sort no, always takes guadratic time

Sorting Challenge 3

Problem: sort a huge number of tiny files (each file is independent) Ex: daily customer transaction records.

Which sorting method to use?

- 1. quicksort
- 2. insertion sort
- 3. selection sort



Sorting Challenge 3

Problem: sort a huge number of tiny files (each file is independent) Ex: daily customer transaction records.

Which sorting method to use?

- 1. quicksort no, too much overhead
- 2. insertion sort YES, much less overhead than system sort
- 3. selection sort YES, much less overhead than system sort

Ex: 4 record file.

- 4 N log N + 35 = 70
- 2N² = 32



Sorting Challenge 4

Problem: sort a huge file that is already almost in order. Ex: re-sort a huge database after a few changes.

Which sorting method to use?

- 1. quicksort probably no, insertion simpler and faster
- 2. insertion sort YES, linear time for most definitions of "in order"

18

20

3. selection sort no, always takes quadratic time

Ex:

- ABCDEFHIJGPKLMNOQRSTUVWXYZ
- ZABCDEFGHIJKLMNOPQRSTUVWXY

Sorting Applications

Sorting algorithms are essential in a broad variety of applications

- Sort a list of names.
- Organize an MP3 library.
- obvious applications • Display Google PageRank results.
- List RSS news items in reverse chronological order.
- Find the median.
- Find the closest pair.
- Binary search in a database.
- Identify statistical outliers.
- Find duplicates in a mailing list.
- Data compression.
- Computer graphics.
- Computational biology.

non-obvious applications

problems become easy once

items are in sorted order

- Supply chain management. • Load balancing on a parallel computer.
- Every system needs (and has) a system sort!

System sort: Which algorithm to use?

Many sorting algorithms to choose from

internal sorts.

- Insertion sort, selection sort, bubblesort, shaker sort.
- Quicksort, mergesort, heapsort, samplesort, shellsort.
- Solitaire sort, red-black sort, splaysort, Dobosiewicz sort, psort, ...

external sorts. Poly-phase mergesort, cascade-merge, oscillating sort.

22

radix sorts.

- Distribution, MSD, LSD.
- 3-way radix guicksort.

parallel sorts.

- Bitonic sort, Batcher even-odd sort.
- Smooth sort, cube sort, column sort,
- GPUsort.

21





Duplicate keys

Often, purpose of sort is to bring records with duplicate keys together.

- Sort population by age.
- Finding collinear points.
- Remove duplicates from mailing list.
- Sort job applicants by college attended.

Typical characteristics of such applications.

- Huge file.
- Small number of key values.

Mergesort with duplicate keys: always ~ N lg N compares

Quicksort with duplicate keys

- algorithm goes quadratic unless partitioning stops on equal keys!
- [many textbook and system implementations have this problem]
- 1990s Unix user found this problem in qsort()

Duplicate keys: the problem

Assume all keys are equal.

Recursive code guarantees that case will predominate!

Mistake: Put all keys equal to the partitioning element on one side

- easy to code
- guarantees N² running time when all keys equal

BAABABCCBCB

A A A A A A A A A

Recommended: Stop scans on keys equal to the partitioning element

- easy to code
- guarantees N lg N compares when all keys equal

Desirable: Put all keys equal to the partitioning element in place

A A A B B B B B C C C A A A A A A A A A A A

Common wisdom to 1990s: not worth adding code to inner loop

3-Way Partitioning

3-way partitioning. Partition elements into 3 parts:

- Elements between \mathtt{i} and \mathtt{j} equal to partition element $\mathtt{v}.$
- No larger elements to left of i.
- No smaller elements to right of j.



Dutch national flag problem.

- not done in practical sorts before mid-1990s.
- new approach discovered when fixing mistake in Unix qsort()
- now incorporated into Java system sort

Solution to Dutch national flag problem.

3-way partitioning (Bentley-McIlroy).

 Partition elements into 4 parts: no larger elements to left of i no smaller elements to right of j equal elements to left of p equal elements to right of q



All the right properties.

- in-place.
- not much code.
- linear if keys are all equal.
- small overhead if no equal keys.





28



```
3-way Quicksort: Java Implementation
          private static void sort(Comparable[] a, int l, int r)
             if (r \le 1) return;
             int i = 1-1, j = r;
             int p = 1-1, q = r;
             while(true)
                                                         4-way partitioning
             -
                while (less(a[++i], a[r])) ;
                while (less(a[r], a[--j])) if (j == 1) break;
                if (i >= j) break;
                exch(a, i, j);
                if (eq(a[i], a[r])) exch(a, ++p, i); swap equal keys to left or right
                if (eq(a[j], a[r])) exch(a, --q, j);
             3
             exch(a, i, r);
                                                         swap equal keys back to middle
             j = i - 1;
             i = i + 1;
             for (int k = 1; k \le p; k++) exch(a, k, j--);
             for (int k = r-1; k \ge q; k--) exch(a, k, i++);
             sort(a, 1, j);
                                                         recursively sort left and right
             sort(a, i, r);
          }
```















Generalized compare

Comparable interface: sort uses type's compareTo () function:

Problem 1: Not type-safe Problem 2: May want to use a different order. Problem 3: Some types may have no "natural" order.

Solution: Use Comparator interface

Comparator interface. Require a method compare () so that compare (v, w) is a total order that behaves like compareto ().

Advantage. Separates the definition of the data type from definition of what it means to compare two objects of that type. • add any number of new orders to a data type.

• add an order to a library data type with no natural order.

Generalized compare

Comparable interface: sort uses type's compareTo() function:

Problem 2: May want to use a different order. Problem 3: Some types may have no "natural" order.

Solution: Use Comparator interface

Example:

37

39

publi {	c class ReverseOrder implements Comparator <string></string>
pu	blic int compare(String a, String b)
ł	return - a.compareTo(b); }
}	reverse sense of comparison
Ar	<pre>crays.sort(a, new ReverseOrder());</pre>
• •	•

Generalized compare Easy modification to support comparators in our sort implementations pass comparator to sort(), less() • use it in less () Example: (insertion sort) public static void sort(Object[] a, Comparator comparator) int N = a.length;for (int i = 0; i < N; i++) for (int j = i; j > 0; j--) if (less(comparator, a[j], a[j-1])) exch(a, j, j-1);else break; } private static boolean less (Comparator c, Object v, Object w) { return c.compare(v, w) < 0; }</pre> private static void exch(Object[] a, int i, int j) { Object t = a[i]; a[i] = a[j]; a[j] = t; }





Generalized compare problem

A typical application

- first, sort by name
- then, sort by section

Arrays.sort(students, Student.BY_NAME);

↓ I						
Andrews	3	Α	664-480-0023	097 Little	Fo	x
Battle	4	С	874-088-1212	121 Whitman	Ch	en
Chen	2	Α	991-878-4944	308 Blair	Kan	aga
Fox	1	Α	884-232-5341	11 Dickinson	Andr	ews
Furia	3	Α	766-093-9873	101 Brown	Fur	ria
Gazsi	4	В	665-303-0266	22 Brown	Roh	ide
Kanaga	3	В	898-122-9643	22 Brown	Bat	tle
Rohde	3	А	232-343-5555	343 Forbes	Ga	zsi

<pre>Arrays.sort(students, Student.BY_SECT);</pre>							
	Ļ						
Fox	1	Α	884-232-5341	11 Dickinson			
Chen	2	Α	991-878-4944	308 Blair			
Kanaga	3	В	898-122-9643	22 Brown			
Andrews	3	Α	664-480-0023	097 Little			
Furia	3	Α	766-093-9873	101 Brown			
Rohde	3	Α	232-343-5555	343 Forbes			
Battle	4	С	874-088-1212	121 Whitman			
Gazsi	4	В	665-303-0266	22 Brown			

42

44

@#%&@!! Students in section 3 no longer in order by name.

A stable sort preserves the relative order of records with equal keys. Is the system sort stable?



Stability

- Q. Which sorts are stable?
- Selection sort?
- Insertion sort?
- Shellsort?
- Quicksort?
- Mergesort?
- A. Careful look at code required.

Annoying fact. Many useful sorting algorithms are unstable.

Easy solutions.

- add an integer rank to the key
- careful implementation of mergesort

Open: Stable, inplace, optimal, practical sort??









- Q. Is the system sort good enough?
- A. Maybe (no matter which algorithm it uses).