

Better than $N \lg N$?

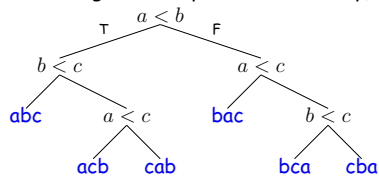
at if all you can do to keys is compare them, then sorting N items takes $\Theta(N \lg N)$ comparisons.

there are $N!$ possible ways the input data could be ordered.

our program must be prepared to do $N!$ different combinations of comparisons and sorting operations.

there must be $N!$ possible combinations of outcomes of comparisons in your program, since those determine what move to make next (we're assuming that comparisons are 2-way).

free time



Beyond Comparison: Distribution

we can do more than compare keys?

if keys are integers, how can we sort a set of N different integer keys whose range is from 0 to kN , for some small constant k ?

the answer is **distribution sorting**:

partition keys into N buckets; integer p goes to bucket $\lfloor p/k \rfloor$. Count keys per bucket, so catenate and use insertion sort, which will now be fast.

Example: $N = 10$:

Input: 10 13 4 2 19 17 0 9
 Counts: 2 | 4 | 1 | 9 | 10 | 13 | 14 | 17 | 19 | 20

Distribution sort is fast. Putting the data in buckets takes time $\Theta(N)$. Insertion sort takes $\Theta(kN)$. When k is fixed (constant), distribution sort runs in time $\Theta(N)$.

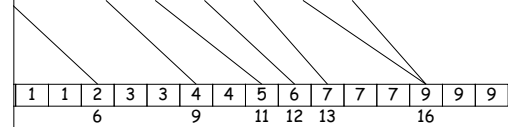
Distribution Counting Example

Items are between 0 and 9 as in this example:

9 1 9 1 9 5 3 7 3 1 6 7 4 2 0

Counts: 2 2 1 1 3 0 3
 3 4 5 6 7 8 9

Running sum: 7 9 11 12 13 16 16
 2 < 3 < 4 < 5 < 6 < 7 < 8 < 9



Counts give # occurrences of each key.

Running sum gives cumulative count of keys < each value...

Running sum tells us where to put each key:

Number of keys < k goes into slot m , where m is the number of keys that are < k ; next k goes into slot $m + 1$, etc.

CS61B Lectures #28

Focus on sorting by comparison

Counting, radix sorts

Today: DS(IJ), Chapter 8; Next topic: Chapter 9.

Necessary Choices

Each k -test goes two ways, number of possible different k -tests is 2^k .

Enough tests so that $2^k \geq N!$, which means $k \geq \lg N!$.

Stirling's approximation,

$$\begin{aligned} N! &\approx \sqrt{2\pi N} \left(\frac{N}{e}\right)^N \left(1 + \Theta\left(\frac{1}{N}\right)\right), \\ \lg N! &\approx \frac{1}{2}(\lg 2\pi + \lg N) + N \lg N - N \lg e + \lg\left(1 + \Theta\left(\frac{1}{N}\right)\right) \\ &= \Theta(N \lg N) \end{aligned}$$

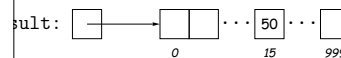
That k , the worst-case number of tests needed to sort by comparison sorting, is in $\Omega(N \lg N)$: there must be cases that need (some multiple of) $N \lg N$ comparisons to sort N items.

Distribution Counting

Unique: count the number of items < 1, < 2, etc.

Items with value < p , then in sorted order, the j^{th} item must be item # $M_p + j$.

If one has a set of numbers in the range $[0, 1000)$ (possibly with duplicates) and that exactly 15 of them are less than 50, which set. Then the result of sorting will look like this:



For each k , the count of numbers < k gives the index of k in the array.

N items in the range $0..M-1$, this gives another **linear-time** algorithm (We include M and N here to allow for cases and for cases where $M \gg N$.)

Notation: the notations $[A, B]$, (A, B) , $[A, B)$, and $(A, B]$ above refer to intervals. The use of square brackets vs. round brackets reflects the distinction between open and closed intervals. $x \in [A, B]$ iff $A \leq x \leq B$, while $x \in (A, B)$ iff $A < x < B$, etc.

Distribution Counting Example (II)

9 1 9 1 9 5 3 7 3 1 6 7 4 2 0

2	2	1	1	3	0	3
3	4	5	6	7	8	9

Counts

7	9	11	12	13	16	16
3	4	5	6	7	8	9

Running sum of Counts

7	9	11	12	14	16	16
3	4	5	6	7	8	9

Next positions

					7									
	6		9		12				15				18	

Output

Distribution Counting Example (II)

9 1 9 1 9 5 3 7 3 1 6 7 4 2 0

2	2	1	1	3	0	3
3	4	5	6	7	8	9

Counts

7	9	11	12	13	16	16
3	4	5	6	7	8	9

Running sum of Counts

7	10	11	12	14	16	16
3	4	5	6	7	8	9

Next positions

			4		7									
	6		9		12				15				18	

Output

Distribution Counting Example (II)

9 1 9 1 9 5 3 7 3 1 6 7 4 2 0

2	2	1	1	3	0	3
3	4	5	6	7	8	9

Counts

7	9	11	12	13	16	16
3	4	5	6	7	8	9

Running sum of Counts

7	10	11	12	14	16	17
3	4	5	6	7	8	9

Next positions

			4		7			9						
	6		9		12			15					18	

Output

Distribution Counting Example (II)

9 1 9 1 9 5 3 7 3 1 6 7 4 2 0

2	2	1	1	3	0	3
3	4	5	6	7	8	9

Counts

7	9	11	12	13	16	16
3	4	5	6	7	8	9

Running sum of Counts

7	9	11	12	13	16	16
3	4	5	6	7	8	9

Next positions

	6		9		12				15				18	

Output

Distribution Counting Example (II)

9 1 9 1 9 5 3 7 3 1 6 7 4 2 0

2	2	1	1	3	0	3
3	4	5	6	7	8	9

Counts

7	9	11	12	13	16	16
3	4	5	6	7	8	9

Running sum of Counts

7	9	11	12	14	16	16
3	4	5	6	7	8	9

Next positions

					7									
	6		9		12				15				18	

Output

Distribution Counting Example (II)

9 1 9 1 9 5 3 7 3 1 6 7 4 2 0

2	2	1	1	3	0	3
3	4	5	6	7	8	9

Counts

7	9	11	12	13	16	16
3	4	5	6	7	8	9

Running sum of Counts

7	10	11	12	14	16	16
3	4	5	6	7	8	9

Next positions

			4		7									
	6		9		12			15					18	

Output

Distribution Counting Example (II)

9 1 9 1 9 5 3 7 3 1 6 7 4 2 0

2	2	1	1	3	0	3
3	4	5	6	7	8	9

Counts

7	9	11	12	13	16	16
3	4	5	6	7	8	9

Running sum of Counts

7	10	11	12	14	16	18
3	4	5	6	7	8	9

Next positions

			4		7		9	9	
	6		9		12		15		18

Output

Distribution Counting Example (II)

9 1 9 1 9 5 3 7 3 1 6 7 4 2 0

2	2	1	1	3	0	3
3	4	5	6	7	8	9

Counts

7	9	11	12	13	16	16
3	4	5	6	7	8	9

Running sum of Counts

7	10	11	12	14	16	19
3	4	5	6	7	8	9

Next positions

1			4		7		9	9	9
	6		9		12		15		18

Output

Distribution Counting Example (II)

9 1 9 1 9 5 3 7 3 1 6 7 4 2 0

2	2	1	1	3	0	3
3	4	5	6	7	8	9

Counts

7	9	11	12	13	16	16
3	4	5	6	7	8	9

Running sum of Counts

8	10	12	12	14	16	19
3	4	5	6	7	8	9

Next positions

1		3		4		5		7		9	9	9
	6		9		12		15		18			

Output

Distribution Counting Example (II)

9 1 9 1 9 5 3 7 3 1 6 7 4 2 0

2	2	1	1	3	0	3
3	4	5	6	7	8	9

Counts

7	9	11	12	13	16	16
3	4	5	6	7	8	9

Running sum of Counts

7	10	11	12	14	16	17
3	4	5	6	7	8	9

Next positions

			4		7		9		
	6		9		12		15		18

Output

Distribution Counting Example (II)

9 1 9 1 9 5 3 7 3 1 6 7 4 2 0

2	2	1	1	3	0	3
3	4	5	6	7	8	9

Counts

7	9	11	12	13	16	16
3	4	5	6	7	8	9

Running sum of Counts

7	10	11	12	14	16	18
3	4	5	6	7	8	9

Next positions

1			4		7		9	9	
	6		9		12		15		18

Output

Distribution Counting Example (II)

9 1 9 1 9 5 3 7 3 1 6 7 4 2 0

2	2	1	1	3	0	3
3	4	5	6	7	8	9

Counts

7	9	11	12	13	16	16
3	4	5	6	7	8	9

Running sum of Counts

7	10	12	12	14	16	19
3	4	5	6	7	8	9

Next positions

1			4		5		7		9	9	9
	6		9		12		15		18		

Output

Distribution Counting Example (II)

9 1 9 1 9 5 3 7 3 1 6 7 4 2 0

2	2	1	1	3	0	3
3	4	5	6	7	8	9

Counts

7	9	11	12	13	16	16
3	4	5	6	7	8	9

Running sum of Counts

9	10	12	12	15	16	19
3	4	5	6	7	8	9

Next positions

1		3	3	4		5		7	7		9	9	9
	6		9		12		15		18				

Output

Distribution Counting Example (II)

9 1 9 1 9 5 3 7 3 1 6 7 4 2 0

2	2	1	1	3	0	3
3	4	5	6	7	8	9

Counts

7	9	11	12	13	16	16
3	4	5	6	7	8	9

Running sum of Counts

9	10	12	13	15	16	19
3	4	5	6	7	8	9

Next positions

1	1		3	3	4		5	6	7	7		9	9	9
	6		9		12		15		18					

Output

Distribution Counting Example (II)

9 1 9 1 9 5 3 7 3 1 6 7 4 2 0

2	2	1	1	3	0	3
3	4	5	6	7	8	9

Counts

7	9	11	12	13	16	16
3	4	5	6	7	8	9

Running sum of Counts

9	11	12	13	16	16	19
3	4	5	6	7	8	9

Next positions

1	1		3	3	4	4	5	6	7	7	7	9	9	9
	6		9		12		15		18					

Output

Distribution Counting Example (II)

9 1 9 1 9 5 3 7 3 1 6 7 4 2 0

2	2	1	1	3	0	3
3	4	5	6	7	8	9

Counts

7	9	11	12	13	16	16
3	4	5	6	7	8	9

Running sum of Counts

8	10	12	12	15	16	19
3	4	5	6	7	8	9

Next positions

1		3		4		5		7	7		9	9	9
	6		9		12		15		18				

Output

Distribution Counting Example (II)

9 1 9 1 9 5 3 7 3 1 6 7 4 2 0

2	2	1	1	3	0	3
3	4	5	6	7	8	9

Counts

7	9	11	12	13	16	16
3	4	5	6	7	8	9

Running sum of Counts

9	10	12	12	15	16	19
3	4	5	6	7	8	9

Next positions

1	1		3	3	4		5		7	7		9	9	9
	6		9		12		15		18					

Output

Distribution Counting Example (II)

9 1 9 1 9 5 3 7 3 1 6 7 4 2 0

2	2	1	1	3	0	3
3	4	5	6	7	8	9

Counts

7	9	11	12	13	16	16
3	4	5	6	7	8	9

Running sum of Counts

9	10	12	13	16	16	19
3	4	5	6	7	8	9

Next positions

1	1		3	3	4		5	6	7	7	7	9	9	9
	6		9		12		15		18					

Output

Distribution Counting Example (II)

9 1 9 1 9 5 3 7 3 1 6 7 4 2 0

2	2	1	1	3	0	3
3	4	5	6	7	8	9

 Counts

7	9	11	12	13	16	16
3	4	5	6	7	8	9

 Running sum of Counts

9	11	12	13	16	16	19
3	4	5	6	7	8	9

 Next positions

1	1	2	3	3	4	4	5	6	7	7	7	9	9	9
		6		9		12		15		18				

 Output

Distribution Counting Example (II)

9 1 9 1 9 5 3 7 3 1 6 7 4 2 0

2	2	1	1	3	0	3
3	4	5	6	7	8	9

 Counts

7	9	11	12	13	16	16
3	4	5	6	7	8	9

 Running sum of Counts

9	11	12	13	16	16	19
3	4	5	6	7	8	9

 Next positions

1	1	2	3	3	4	4	5	6	7	7	7	9	9	9
		6		9		12		15		18				

 Output

MSD Radix Sort

Complicated: must keep lists from each step separate
 processing 1-element lists

A	posn
cat, cad, con, bat, can, be, let, bet	0
be, bet / cat, cad, con, can / let / set	1
* be, bet / cat, cad, con, can / let / set	2
be / bet / * cat, cad, con, can / let / set	1
be / bet / * cat, cad, can / con / let / set	2
be / bet / cad / can / cat / con / let / set	

divide partially sorted sublists that will never be moved
 be occupied by other sublists.

mark a sublist to be sorted on some character position.

And Don't Forget Search Trees

A search tree is in sorted order, when read in inorder.

use B-trees to really use for sorting [next topic].

Same performance as heapsort: N insertions in time
 is $\Theta(N)$ to traverse, gives

$$\Theta(N + N \lg N) = \Theta(N \lg N)$$

Radix Sort

Sorts *one character at a time*.

Distribution counting for each digit.

Order from right to left (LSD radix sort) or left to right (MSD)

Radix sort is venerable: used for punched cards. Example:

Initial: set, cat, cad, con, bat, can, be, let, bet

		bet				bat	bet
		let				cat	let
		bat				can	set
		cat				cad	be
		set				con	con
e	cad	con	set				
'	'd'	'n'	't'			'd'	'e'

Pass 2 (by char #1)

h, can, set, cat, bat, let, bet cad, can, cat, bat, be, set, let, bet, con

		con					
		bet	cat				
		be	can				
		bat	cad	let	set		
		'b'	'c'	'l'	's'		

Pass 3 (by char #0)

bat, be, bet, cad, can, cat, con, let, set

Performance of Radix Sort

Sorts takes $\Theta(B)$ time where B is *total size of the key data*.

Compare other sort times as functions of *#records*.

Are there?

For different records, one must have keys at least $\Theta(\lg N)$

Radix sort, comparison actually takes time $\Theta(K)$ where K is size
 of the largest case [why?]

Radix sort comparisons really means $N(\lg N)^2$ operations.

Radix sort would take $B = N \lg N$ time for N records with
 $\Theta(\lg N)$ keys.

On the other hand, we must work to get good constant factors with

Summary

Insertion sort: $\Theta(Nk)$ comparisons and moves, where k is maximum displacement of an element from its final position.

Works well on small datasets or almost ordered data sets.

Heap sort: $\Theta(N \lg N)$ with good constant factor if data is not pathological. $\mathcal{O}(N^2)$.

Merge sort: $\Theta(N \lg N)$ guaranteed. Good for external sorting.

Quick sort with guaranteed balance: $\Theta(N \lg N)$ guaranteed.

Distribution sort: $\Theta(B)$ (number of bytes). Also good for external sorting.