

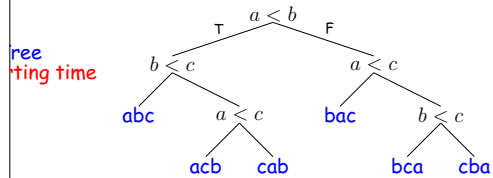
Better than $N \lg N$?

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

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

our program must be prepared to do $N!$ different comparisons and data-moving 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).



Beyond Comparison: Distribution

can we do more than compare keys?

how can we sort a set of N integer keys whose values are in the range $[0, kN]$, for some small constant k ?

idea: put the integers into N buckets, with an integer p per bucket, where $p \leq [p/k]$.

sort keys per bucket, so concatenate and use insertion sort, which is fast.

Example: $p = 10$:

10 13 4 2 19 17 0 9
 buckets: 2 | 4 | 9 | 10 | 13 | 14 | 17 | 19 |

bucket sort is fast. Putting in buckets takes time $\Theta(N)$, and sorting buckets takes $\Theta(kN)$. When k is fixed (constant), we have total 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
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

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

Counts

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

Running sum

1	1	2	3	3	4	4	5	6	7	7	7	9	9	9
		6		9		11	12	13				16		

counts gives # occurrences of each key.

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

running sum tells us where to put each key:

the number of keys k goes into slot m , where m is the number of keys that are < k .

CS61B Lectures #28

focus on sorting by comparison

distribution counting, radix sorts

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

Necessary Choices

if we use k -tests goes two ways, number of possible different outcomes of k -tests is 2^k .

enough tests so that $2^k \geq N!$, which means $k \in \Omega(\lg N!)$.

Stirling's approximation,

$$\sqrt{2\pi N} \left(\frac{N}{e}\right)^N \left(1 + \Theta\left(\frac{1}{N}\right)\right),$$

$$\frac{1}{2}(\lg 2\pi + \lg N) + N \lg N - N \lg e + \lg \left(1 + \Theta\left(\frac{1}{N}\right)\right)$$

$$\Theta(N \lg N)$$

that k , the worst-case number of tests needed to sort N items by comparison sorting, is in $\Omega(N \lg N)$: there must be cases that require (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 # $M_p + j$.

for linear-time algorithm.

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

						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	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	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	

And Don't Forget Search Trees

tree is in sorted order, when read in inorder.

useful to really use for sorting [next topic].

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

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

Radix Sort

processes one character at a time.

distribution counting for each digit.

order right to left (LSD radix sort) or left to right (MSD)

Radix sort is venerable: used for punched cards.

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

	bet			bat	bet
	let			cat	let
	bat			can	set
	cat			cad	be
	con			con	con
Pass 2					
(by char #1)					
set	cad	con	set	cad	be
'd'	'n'	't'		'a'	'e'
				'o'	
can, set, cat, bat, let, bet				cad, can, cat, bat, be, set, let, bet, con	
Pass 3					
(by char #0)					
	bet	con			
	be	cat			
	bat	can			
	cad	let	set		
	'b'	'c'	'l'	's'	
	bat, be, bet, cad, can, cat, con, let, set				

Performance of Radix Sort

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

Compare other sorts as function of #records.

Are there any?

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

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

Number of comparisons really means $N(\lg N)^2$ operations.

Radix sort would take $B = N \lg N$ time with minimal-length

On the other hand, 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. Worst case $O(N^2)$.

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.