how fast can we sort ?
it depends on what we call the computational model
      ‘what you are allowed to do with the elements’

can we do better than $\Theta(n\lg n)$ ?
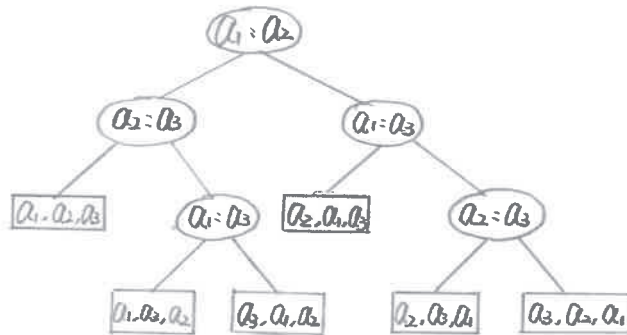
Comparison Sorting Model:
      only use comparisons to determine the relative order of elements

Decision - Tree Model: (决策树)
      more general than comparison model
      Example:
            to sort  $<a_1, a_2, a_3>$

$$a_1:a_2$$
$$a_2:a_3 \qquad a_1:a_3$$
$$\boxed{a_1,a_2,a_3} \quad a_1:a_3 \quad \boxed{a_2,a_1,a_3} \quad a_2:a_3$$
$$\boxed{a_1,a_3,a_2} \quad \boxed{a_3,a_1,a_2} \qquad \boxed{a_2,a_3,a_1} \quad \boxed{a_3,a_2,a_1}$$

      what this tree means is that each node you're making a comparison $(x:y)$.
      if $x<y$, go left, and go right otherwise. when you get down to
      a leaf, this is the answer.

   Def. in general, $<a_1, a_2, \cdots, a_n>$,
      △ each internal node ( non-leaf node) has a label of the form $"i:j"$ where
         $1 \leq i, j \leq n$, means we compare $a_i$ and $a_j$,
         and we have two subtrees from every such node,
               • left subtree which tells you what the algorithm does when $a_i \leq a_j$
               • right subtree which gives subsequent comparisons if $a_i > a_j$
      △ each leaf node gives a permutation, $<\pi(1), \pi(2), \cdots, \pi(n)>$,
         such that  $a_{\pi(1)} \leq a_{\pi(2)} \leq \cdots \leq a_{\pi(n)}$

用决策树的表达方式 构建比较排序算法。(转换的过程)      “基于比较的排序可以转换成决策树。
                                                    但是有一些排序算法无法以决策树的形式
                                                    表现出来”
• one tree for each $n$
• view algorithm as splitting into two forks (subtrees) whenever it makes a comparison
• tree lists comparisons along all possible instruction traces
• running time (the number of comparisons) = the length of path
• the worst-case running time = the height of the tree

<u>lower bound on decision-tree sorting:</u>

"定理" any decision-tree sorting $n$ elements has height $\underline{\Omega(n \lg n)}$ "$\geq$"

    proof: the number of leaves must be at least $n!$
        (as there are $n$ factorial permutations of an input)
        the height of the tree $:= h$, then it has at most $2^h$ leaves

$$\Rightarrow \#\text{leaves} \leq 2^h$$
$$\Rightarrow n! \leq 2^h$$
$$\Rightarrow h \geq \log_2(n!)$$
$$\text{"斯特林公式"} \geq \log_2\left(\frac{n}{e}\right)^n$$
$$= n \log_2 \frac{n}{e}$$
$$= n \log_2 n - n \log_2 e$$
$$= \Omega(n \lg n)$$

corollary: merge sort and heapsort are asymptotically optimal $(n \lg n)$,
    but this is only in the comparison model , "基于比较模型的话，
    Randomized QuickSort is too in expectation. $n \lg n$ 就是极限了"

Sorting in Linear Time "不可能比线性时间更快完成排序，因为得遍历数据"
$\geq$ algorithms for instance

① counting sort
    input : $A[1 \ldots n]$ , each $A[i]$ is an integer from the range of 1 to $k$
    output : $B[1 \ldots n]$ = sorting of $A$ "当k较小时，性能比较好"
    auxiliary storage : $C[1 \ldots k]$

    Counting Sort:

```
for i ← 1 to k
    do C[i] ← 0
for j ← 1 to n
    do C[A[j]] ← C[A[j]]+1       // C[i]表示数值i出现的次数
for i ← 2 to k
    do C[i] ← C[i] + C[i-1]       // C[i]表示 小于等于i的元素的数目（对前缀做加法）
for j ← n downto 1                                          "prefix sum"
    do B[C[A[j]]] ← A[j]          // distribution
        C[A[j]] ← C[A[j]]-1
```

    Example:
$$A = [4, 1, 3, 4, 3]$$

$C =$ | 0 | 0 | 0 | 0 |  $\xrightarrow{\text{counting}}$  $C =$ | 1 | 0 | 2 | 2 |
       1   2   3   4                                        1   2   3   4

$\xrightarrow{\text{prefix sum}}$ $C' =$ | 1 | 1 | 3 | 5 |
                                          1   2   3   4

$j=5$ , $A[j]=3$ , $B =$ | | | 3 | | |
                              $C'[A[j]]=3$ $\longleftarrow$ "值为3的元素应该放在位置3或者更靠前"
                              $C' = [1, 1, ②, 5]$

$j=4$ , $A[j]=4$ , $B =$ | | | 3 | 4 |
                              $C'[A[j]]=5$
                              $C' = [1, 1, 2, ④]$

$\cdots \cdots \cdots$

$B =$ | 1 | 3 | 3 | 4 | 4 |

$T(n) = O(k+n)$ ,
it is a great algorithm if $k$ is relatively small , like at most $n$ .
so you could write a combination algorithm that if $k > n\lg n$ and if $k \leq n\lg n$
a stable sorting algorithm preserves the ~~the~~ relative order of equal elements , counting sort is stable

② radix sort 基数排序 by Herman Hollerith in 1890
radix sort ~~is going to~~ work for a much larger range of numbers in linear time
first: sort by the most significant digit first (need many boxes)
right: (by Hollerith) sort by the least significant digit first, using stable sorting

> Hollerith 在1911年创建制表机公司 (tabulating machine company),
> 然后在1924年合并了其他几个公司,组成了 IBM

the whole idea is that we are doing a digit-by-digit sort ,
from least significant digit to most significant digit
the nice thing about this algorithm is that there are no bins , it's one big bin at all times
Example:

```
3 2 9        3 2 9              7 2 0
4 5 7        4 5 7              3 5 5
6 5 7        6 5 7              4 3 6
8 3 9  "LSB" 8 3 9   "stable"   4 5 7
4 3 6   -->  4 3 6     -->      6 5 7
7 2 0        7 2 0              3 2 9
3 5 5        3 5 5              8 3 9
```

$$
\text{"stable"} \rightarrow
\begin{array}{ccc}
7 & 2 & 0 \\
3 & 2 & 9 \\
4 & 3 & 6 \\
8 & 3 & 9 \\
3 & 5 & 5 \\
4 & 5 & 7 \\
6 & 5 & 7 \\
\end{array}
\quad \text{"stable"} \rightarrow
\begin{array}{ccc}
3 & 2 & 9 \\
3 & 5 & 5 \\
4 & 3 & 6 \\
4 & 5 & 7 \\
6 & 5 & 7 \\
7 & 2 & 0 \\
8 & 3 & 9 \\
\end{array} \checkmark
$$

("sorted")

"when I have equal elements here, I have already sorted the suffix"

"好的部分是我们不用分成一个一个箱子了，而是始终把它们放在一个大箱子里面。"

Correctness:
to induct on the digit position $t$ that we are currently sorting,
assume that by induction that it is already sorted on lower $t$-1 digits,
and then the next thing we do is to sort on the $t$-th digit.
if two elements are the same (has the same $t$-th digit),
stability $\Rightarrow$ keep the order $\Rightarrow$ still sorted
else,
put them in the right order

next we are going to use (counting sort) for each round
(we could use any sorting algorithm we want for individual digits)

Analysis:
- use counting sort ($O(k+n)$)

"可以把几个比特放在一起当做一位数处理，相当于八进制、十六进制"

- say $n$ integers, each $b$ bits long ($0 \sim 2^b - 1$)
- split each integer into $b/r$ "digits", each "digit" is $r$ bits long

"需要运算的轮数"

"基于 $2^r$ 进制来表示这个数"
"counting sort 的 $k = 2^r$"

$$T(n) = O\left(\frac{b}{r} \cdot (n + 2^r)\right)$$

$\min\limits_{r} T(n)$

$r = \log_2 n$, $T(n) = O\left(\frac{bn}{\log_2 n}\right)$

if numbers (integers) are in the range $0 \sim 2^b - 1$
$0 \sim n^d - 1$
then $T(n) = O(d \cdot n)$