

递归与分治

陈长建

计算机科学系

回顾 - 大整数乘法

- 请设计一个有效的算法，可以进行两个n位大整数的乘法运算

— 分治法:

$$X = \begin{array}{|c|c|} \hline \text{n/2位} & \text{n/2位} \\ \hline a & b \\ \hline \end{array} \quad Y = \begin{array}{|c|c|} \hline \text{n/2位} & \text{n/2位} \\ \hline c & d \\ \hline \end{array}$$

— $XY = (a*c) \cdot 10^n + (a*d + b*c) \cdot 10^{n/2} + b*d$

— $XY = a*c \cdot 10^n + ((a-c)*(b-d) + a*c + b*d) \cdot 10^{n/2} + b*d$

回顾 - 大整数乘法

- 更快的方法：
 - 如果将大整数分成更多段，用更复杂的方式把它们组合起来，将有可能得到更优的算法。
 - 最终，这个思想导致了快速傅利叶变换(Fast Fourier Transform)的产生。该方法也可以看作是一个复杂的分治算法，对于大整数乘法，它能在 $O(n \log n)$ 时间内解决。
- 是否能找到线性时间算法？ 目前为止还没有结果

回顾-分治法时间复杂度分析

- 递归树方法
 - 二分搜索算法
- 主定理法

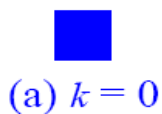
$$T(n) = \begin{cases} c & n = 1 \\ aT\left(\frac{n}{b}\right) + cn & n > 1 \end{cases}$$

$$T(n) = \begin{cases} \Theta(n) & a < b \\ \Theta(n \log_b n) & a = b \\ \Theta(n^{\log_b a}) & a > b \end{cases}$$

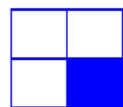
- 代换法 (数学归纳)

回顾 - 棋盘覆盖

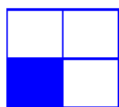
- 在一个 $2^k \times 2^k$ 个方格组成的棋盘中，恰有一个方格与其他方格不同，称该方格为一特殊方格，且称该棋盘为一特殊棋盘。在棋盘覆盖问题中，要用图示的4种不同形态的L型骨牌覆盖给定的特殊棋盘上除特殊方格以外的所有方格，且任何2个L型骨牌不得重叠覆盖。



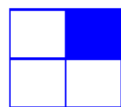
(a) $k = 0$



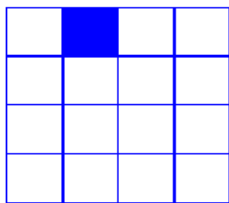
(b) $k = 1$



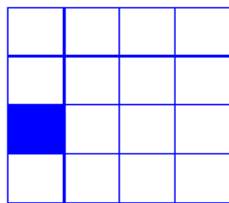
(c) $k = 1$



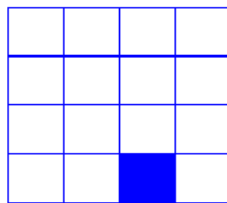
(d) $k = 1$



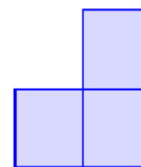
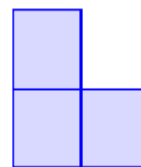
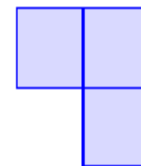
(e) $k = 2$



(f) $k = 2$

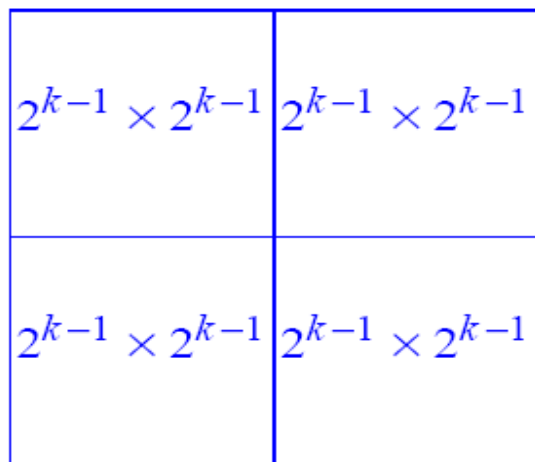


(g) $k = 2$

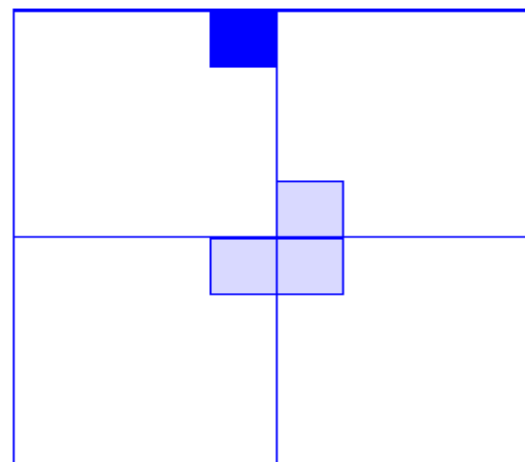


回顾 - 棋盘覆盖

- 当 $k > 0$ 时, 将 $2^k \times 2^k$ 棋盘分割为4个 $2^{k-1} \times 2^{k-1}$ 子棋盘(a)所示。特殊方格必位于4个较小子棋盘之一中, 其余3个子棋盘中无特殊方格。为了将这3个无特殊方格的子棋盘转化为特殊棋盘, 可以用一个L型骨牌覆盖这3个较小棋盘的会合处, 如 (b)所示, 从而将原问题转化为4个较小规模的棋盘覆盖问题。递归地使用这种分割, 直至棋盘简化为棋盘 1×1 。



(a) Partitioning



(b) Triomino placed

时间复杂度分析

$$T(k) = \begin{cases} O(1) & k = 0 \\ 4T(k-1) + O(1) & k > 0 \end{cases}$$

设M(k)为chessBoard算法在计算覆盖一个 $2^k \times 2^k$ 棋盘所需时间:

当 $k > 1$ 时, $M(k) = 4M(k-1)$, $M(0) = 1$

$$M(k) = 4M(k-1) \quad \text{替换 } M(k-1) = 4M(k-2)$$

$$= 4[4M(k-2)] = 4^2 M(k-2)$$

=

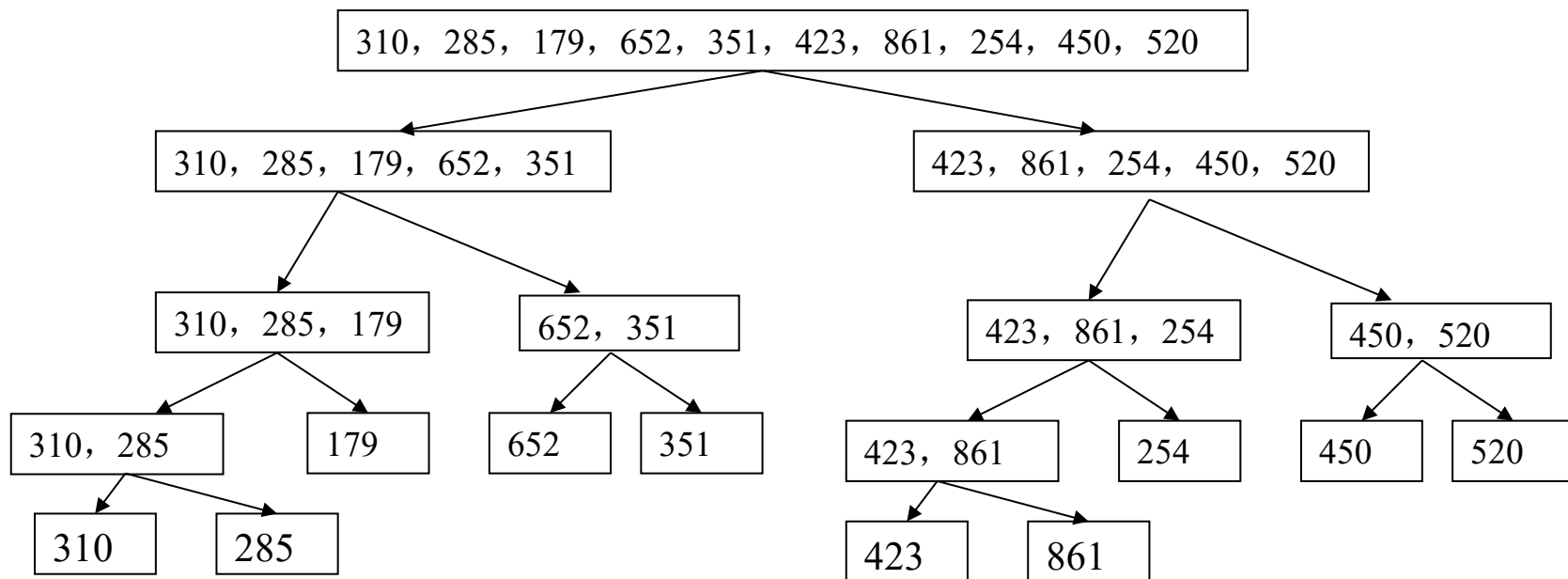
$$= 4^k M(k-k) = 4^k$$

$$T(K) = O(4^K)$$

合并排序基本思想

- 将待排序元素分成大小大致相同的2个子集合，分别对2个子集合进行排序，最终将排好序的子集合合并成为所要求的排好序的集合。

— 例：A = (310, 285, 179, 652, 351, 423, 861, 254, 450, 520)



性能分析

- 1) 空间特性
 - $2n+3$ 个空间位置
- 2) 算法merge的时间与两数组元素的总数成正比
 - (可表示为: cn , n 为元素个数, c 为某正常数)
- 3) 算法mergeSort的时间用递推关系式表示如下:

$$T(n) = \begin{cases} a & n=1, a \text{是常数} \\ 2T(n/2) + cn & n>1, c \text{是常数} \end{cases}$$

以比较为基础分类的时间下界

任何以关键字比较为基础的分类算法，其最坏情况下的时间下界都为： $\Omega(n \log n)$

利用二元比较树证明。

假设参加分类的 n 个关键字 $A(1), A(2), \dots, A(n)$ 互异。任意两个关键字的比较必导致 $A(i) < A(j)$ 或 $A(i) > A(j)$ 的结果。

以二元比较树描述元素间的比较过程：

- 若 $A(i) < A(j)$ ，进入下一级的左分支
- 若 $A(i) > A(j)$ ，进入下一级的右分支

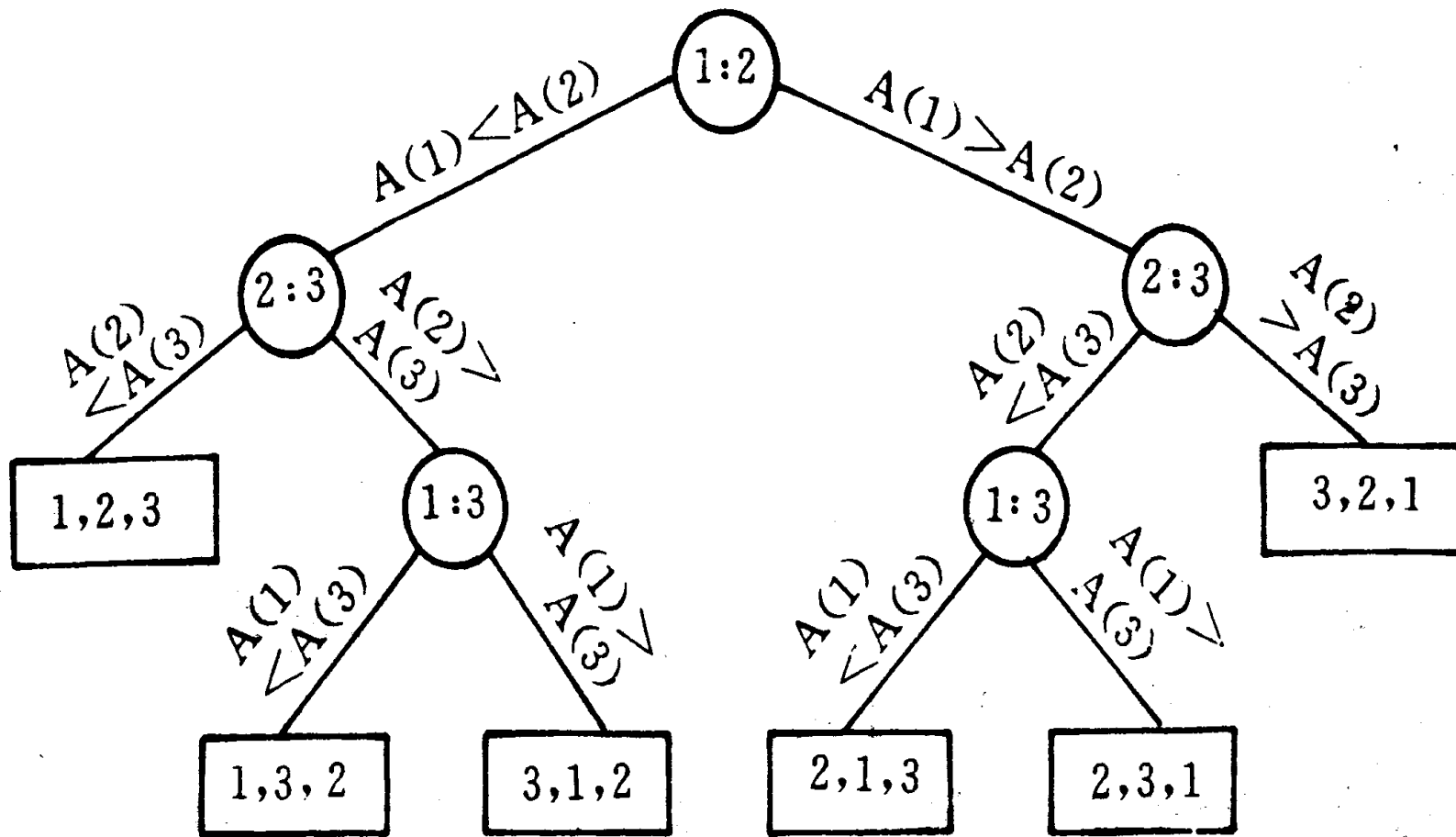


图 2.6 对三个关键字分类的比较树

算法在**外部结点**终止。

从根到某外结点的**路径**代表某个特定输入情况
下一种唯一的**分类排序序列**。路径长度表示生成该
序列代表的分类表所需要的**比较次数**。而**最长的路
径**代表算法在最坏情况下的执行情况，该路径的长
度即是算法在最坏情况下所作的比较次数。

故，以比较为基础的分类算法的最坏情况下界等
于该算法对应的比较树的**最小高度**。

① 由于 n 个关键字有 $n!$ 种可能的排列，所以二元比较树中将有 $n!$ 个外部结点：每种排列对应于某种特定输入情况下的分类情况，每个外部结点表示一种可能的分类序列。

② 设一棵二元比较树的所有内结点的级数均小于或等于 k ，则该树中最多有 2^k 个外结点。

记算法在最坏情况下所作的比较次数为 $T(n)$ ，则有 $T(n)=k$ ：生成外结点所代表的分类序列所需的比较次数等于该外结点所在的级数-1；

根据①和②的分析，有： $n! \leq 2^{T(n)}$

化简：

当 $n > 1$ 时，有 $n! \geq n(n-1)(n-2) \cdots (\lceil n/2 \rceil) \geq (n/2)^{n/2}$

当 $n \geq 4$ 时，有 $T(n) \geq (n/2) \log(n/2) \geq (n/4) \log n$

故，任何以比较为基础的分类算法的最坏情况的时间下界为：

$$\Omega(n \log n)$$

$$n! \approx \sqrt{2\pi n} \left(\frac{n}{e}\right)^n \leq 2^{T(n)}$$

本章要点

- 算法总体思想
- 递归的概念
- 递归案例分析
- 分治法的概念
- 分治案例分析
 - 二分搜索、大整数法、棋盘覆盖、合并排序、快速排序、线性时间选择 ...

分治-例5 快速排序

- 任何一个基于比较来确定两个元素相对位置的排序算法需要 $\Omega(n \log n)$ 计算时间。如果我们能设计一个需要 $O(n \log n)$ 时间的排序算法，则在渐近的意义下，这个排序算法就是最优的。许多排序算法都是追求这个目标。
- 下面介绍快速排序算法，它在平均情况下需要 $(n \log n)$ 时间。这个算法是于1962年由C.A.R.Hoare提出的。

基本思想

- 快速排序是一种基于划分的排序方法;
- 划分: 选取待分类集合A中的某个元素 t , 按照与 t 的大小关系重新整理A中元素, 使得整理后的序列中所有在 t 以前出现的元素均小于等于 t , 而所有出现在 t 以后的元素均大于等于 t 。这一元素的整理过程称为划分 (Partitioning)。元素 t 称为划分元素。
- 快速排序: 通过反复地对待排序集合进行划分达到排序目的的排序算法。

划分过程的算法描述

用元素 $a[p]$ 划分集合 $a[p : r]$

```
public static int partition (int p, int r)
{
    int i = p, j = r + 1;
    Comparable x = a[p];
    // 将 $\geq x$ 的元素交换到左边区域
    // 将 $\leq x$ 的元素交换到右边区域
    while (true) {
        while (a[++i].compareTo(x) < 0); // i由左向右移
        while (a[--j].compareTo(x) > 0); // j由右向左移
        if (i >= j) break;
        MyMath.swap(a, i, j);
    }
    a[p] = a[j];
    a[j] = x;
    return j; //划分元素在位置j
}
```

例子 划分实例

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	i	j
A:	65	70	75	80	85	60	55	50	45	$+\infty$	2	9

|.....|

A:	65	45	75	80	85	60	55	50	70	$+\infty$	3	8
----	----	----	----	----	----	----	----	----	----	-----------	---	---

|.....|

A:	65	45	50	80	85	60	55	75	70	$+\infty$	4	7
----	----	----	----	----	----	----	----	----	----	-----------	---	---

|.....|

A:	65	45	50	55	85	60	80	75	70	$+\infty$	5	6
----	----	----	----	----	----	----	----	----	----	-----------	---	---

|.....|

A:	65	45	50	55	60	85	80	75	70	$+\infty$	6	5
----	----	----	----	----	----	----	----	----	----	-----------	---	---

交换划分元素

→ |.....|

A:	60	45	50	55	65	85	80	75	70	$+\infty$		
----	----	----	----	----	----	----	----	----	----	-----------	--	--

↑
划分元素定位于此

分析

- 经过一次“划分”后，实现了对集合元素的调整：其中一个子集合的所有元素均小于等于另外一个子集合的所有元素。
- 按同样的策略对两个子集合进行分类处理。当子集合分类完毕后，整个集合的分类也完成了。这一过程避免了子集合的归并操作。这一分类过程称为快速分类。

快速排序算法实现

通过反复使用划分算法 **partition** 实现对集合元素的排序。

以 $a[p]$ 为基准元素将 $a[p:r]$ 划分成3段：

$a[p:q-1]$, $a[q]$, $a[q+1:r]$,

使得： $a[p:q-1]$ 中任何元素小于等于 $a[q]$,

$a[q+1:r]$ 中任何元素大于等于 $a[q]$,

下标 q 在划分过程中确定。

```
public static void qSort(int p, int r)
{
    if (p < r) {
        int q = partition(p, r);
        qSort (p, q-1); //对左半段排序
        qSort (q+1, r); //对右半段排序
    }
}
```

快速排序分析

- 记录比较和交换是从两端向中间进行
 - 关键字较大的记录一次就能交换到后面单元
 - 关键字较小的记录一次就能交换到前面单元
 - 记录每次移动的距离较大，因而总比较和移动次数较少。
- 快速排序算法的**性能**取决于划分的**对称性**。
 - 可以设计出采用随机选择策略的快速排序算法。
 - 在快速排序算法的每一步中，当数组还没有被划分时，可以在 $a[p:r]$ 中随机选出一个元素作为划分基准
 - 这样可以使划分基准的选择是随机的，从而可以期望划分是较对称的。

中值快速排序

- 中值快速排序 (median-of-three quick sort) 是快速排序一种变化, 这种算法有更好的平均性能。在快速排序中总是选择 $a[1]$ 做为支点, 而在中值快速排序算法中, 取 $\{a[1], a[(1+r)/2], a[r]\}$ 中大小居中的那个元素作为支点。
- 实现中值快速排序算法的一种最简单的方式就是首先选出中值元素并与 $a[1]$ 进行交换, 然后利用快速排序算法完成排序。

时间复杂度分析

- 统计的对象：元素的比较次数，记为： $C(n)$
- 两点假设
 - ①参加分类的 n 个元素各不相同
 - ② **partition** 中的划分元素 t 是随机选取的（针对平均情况的分析）

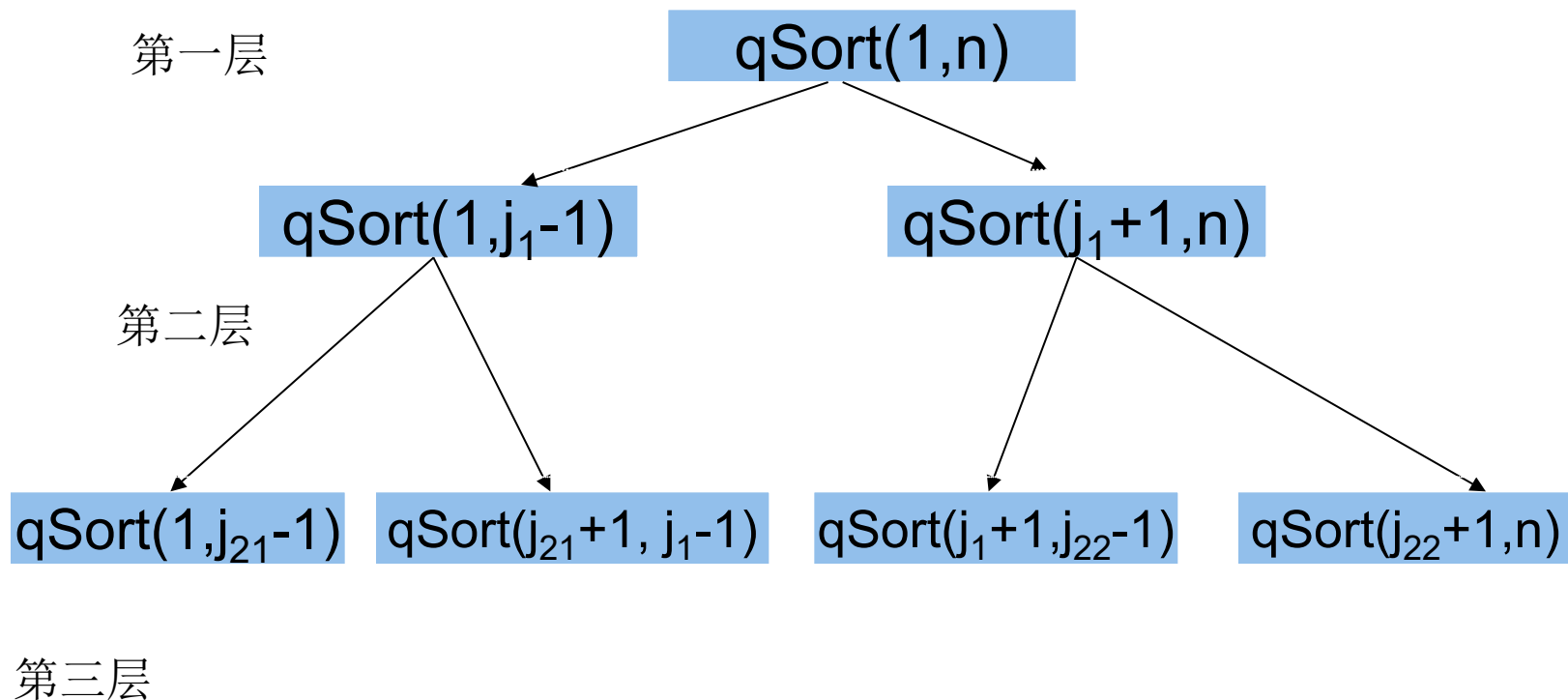
➤ 随机选取划分元素：

在划分区间 $[p, r]$ 随机生成某一坐标： $i \leftarrow \text{random}(p, r)$;

调换 $a[p]$ 与 $a[i]$

- 作用：将随机指定的划分元素的值依旧调换到 $a[p]$ 位置。之后，算法主体不变，仍从 $a[p]$ 开始执行划分操作。

递归层次



设在任一级递归调用上，调用**partition**处理的所有元素总数为 r ，则，初始时 $r=n$ ，以后的每级递归上，由于删去了上一级的划分元素，故 r 比上一级至少1：
理想情况，第一级少1，第二级少2，第三级少4， ...；
最坏情况，每次仅减少1（如集合元素已经按照递增或递减顺序排列）

最坏情况分析

- 记最坏情况下的元素比较次数是 $C_w(n)$;
- **partition**一次调用中的元素比较数是 $r - p + 1$, 故每级递归调用上, 元素的比较次数等于该级所处理的待分类元素个数。
- 最坏情况下, 每级递归调用的元素总数仅比上一级少1, 故 $C_w(n)$ 是 r 由 n 到2的累加和。

即:

$$C_w(n) = \sum_{2 \leq i \leq n} i = O(n^2)$$

最好情况分析

- 记最好情况下的元素比较次数是 $C_{best}(n)$;
- 最好情况下, 每级递归调用的元素总数仅为上一级的1/2, 故 $C_{best}(n)$ 满足递推式:

$$C_{best}(n) = \begin{cases} 0 & n = 1 \\ 2C_{best}(n/2) + n & n > 1 \end{cases}$$

$$\begin{aligned} n > 1, \quad C_{best}(n) &= 2C_{best}(n/2) + n \\ &= 2(2C_{best}(n/2^2) + n/2) + n \\ &= 2^2 C_{best}(n/2^2) + 2n = \dots \\ &= 2^k C_{best}(1) + kn \end{aligned}$$

$$\text{当 } n = 2^k, C_{best}(n) = n \log n$$

平均情况分析

- 记平均情况下的元素比较次数是 $C_A(n)$;
 - 平均情况是指集合中的元素以任意一种顺序排列, 且任选所有可能的元素作为划分元素进行划分和分类, 在这些所有可能的情况下, 算法执行性能的平均值。
 - 设调用**partition**(p,r)时, 所选取划分元素t 恰好是 $a[p:r]$ 中的第i小元素 ($1 \leq i \leq r-p$) 的概率相等。则经过一次划分, 所留下的待分类的两个子数组恰好是 $a[p:j-1]$ 和 $a[j+1:p-1]$ 的概率是: $1/(r-p)$, $p \leq j < r$ 。则有,

$$C_A(n) = n + 1 + \frac{1}{n} \sum_{1 \leq k \leq n} (C_A(k-1) + C_A(n-k))$$

- $n+1$ 是**partition**第一次调用时所需的元素比较次数。
- $C_A(0) = C_A(1) = 0$

$$C_A(n) = n + 1 + \frac{1}{n} \sum_{1 \leq k \leq n} (C_A(k-1) + C_A(n-k))$$



平均情况分析

- 化简上式可得:

$$\begin{aligned} C_A(n)/(n+1) &= C_A(n-1)/n + 2/(n+1) \\ &= C_A(n-2)/(n-1) + 2/n + 2/(n+1) \\ &= C_A(n-3)/(n-2) + 2/(n-1) + 2/n + 2/(n+1) \\ &\dots \\ &= C_A(1)/2 + 2 \sum_{3 \leq k \leq n+1} 1/k \end{aligned}$$

由于 $\sum_{3 \leq k \leq n+1} 1/k \leq \int_2^{n+1} \frac{dx}{x} < \log_e(n+1)$

所以得, $C_A(n) < 2(n+1)\log_e(n+1) = O(n \log n)$

各种排序算法的比较

方法	最坏复杂性	平均复杂性
冒泡排序	n^2	n^2
计数排序	n^2	n^2
插入排序	n^2	n^2
选择排序	n^2	n^2
堆排序	$n \log n$	$n \log n$
归并排序	$n \log n$	$n \log n$
快速排序	n^2	$n \log n$

分治-例6 线性时间选择

- 问题描述

- 给出含有 n 个元素表 $A[0:n-1]$ ，要求确定其中的第 k 小元素。

- 设计思路

- 选择问题可在 $O(n \log n)$ 时间内解决，方法是首先对这 n 个元素进行排序(如使用堆排序或归并排序)，然后取出 $A[k-1]$ 中的元素。若使用快速排序，可以获得更好的平均性能。尽管该算法有一个比较差的渐近复杂性 $O(n^2)$ 。

线性时间选择

- 利用 **partition** 函数。如果划分元素 v 测定在 $A(j)$ 的位置上, 则有 $j-1$ 个元素小于或等于 $A(j)$, 且有 $n-j$ 个元素大于或等于 $A(j)$ 。此时,
 - 若 $k=j$, 则 $A(j)$ 即是第 k 小元素; 否则,
 - 若 $k < j$, 则第 k 小元素将出现在 $A(0:j-1)$ 中;
 - 若 $k > j$, 则第 k 小元素将出现在 $A(j+1:n-1)$ 中。

算法描述

- public static void **randomizedSelect** (int p, int r,int k)
{
 if(p==r) return a[p];
 int i = **randomizedPartition**(p,r),
 j=i-p+1;
 if(k<=j) return **randomizedSelect** (p, i, k); //在左半段
 else return **randomizedSelect** (i+1, r, k-j); //在右半段
}

算法分析

- 两点假设

- ① A中的元素互异

- ② 随机选取划分元素，且选择A中任一元素作为划分元素的概率相同

- 分析

- **randomizedPartition(p,r)**,所需的元素比较次数是 $O(r-p+1)$ 。

- 在执行一次**randomizedPartition(p,r)**后，或者找到第k小元素，或者将在缩小的子集 $A(p, i)$ 或 $A(i+1, r)$ 中继续查找。缩小的子集的元素数将至少比上一次划分的元素数少1。

1) 最坏情况

randomizedSelect的最坏情况时间是 $O(n^2)$

当A中的元素已经按照**递增**的顺序排列, 且 **$k=n$**
此时, 需要 **n 次**调用**randomizedPartition**过程, 且
每次返回的元素位置是子集中的第一个元素, 子集
合的元素数一次**仅减少 1**

则, n 次调用的时间总量是

$$O\left(\sum_{1}^n (i+1)\right) = O(n^2)$$

2) 平均情况

设 $T_A^k(n)$ 是找 $A(1:n)$ 中第 k 小元素的平均时间。 $T_A(n)$ 是 **randomizedSelect** 的平均计算时间, 则有

$$T_A(n) = \frac{1}{n} \sum_{1 \leq k \leq n} T_A^k(n)$$

并定义 $R(n) = \max_k \{T_A^k(n)\}$

则有: $T_A(n) \leq R(n)$ 。

定理2.4 *randomizedSelect*的平均计算时间 $T_A(n)$ 是 $O(n)$

证明:

执行**randomizedPartition**时, 比较语句的执行时间是 $O(n)$ 。
在随机等概率选择划分元素时, 首次调用**randomizedPartition**
中划分元素 v 刚好是 A 中第 i 小元素的概率为 $1/n$, $1 \leq i \leq n$ 。

则, 存在正常数 c , $c > 0$, 有,

$$T_A^k(n) \leq cn + \frac{1}{n} \left(\sum_{1 \leq i < k} T_A^{k-i}(n-i) + \sum_{k < i \leq n} T_A^k(i-1) \right) \quad n \geq 2$$

$$\begin{aligned} \text{且有, } R(n) &\leq cn + \frac{1}{n} \max_k \left\{ \sum_{1 \leq i < k} R(n-i) + \sum_{k < i \leq n} R(i-1) \right\} \\ &= cn + \frac{1}{n} \max_k \left\{ \sum_{n-k+1}^{n-1} R(i) + \sum_k^{n-1} R(i) \right\} \quad n \geq 2 \end{aligned}$$

令 $c \geq R(1)$ 。利用**数学归纳法**证明, 对所有 $n \geq 2$, 有 $R(n) \leq 4cn$ 。

① 当 $n=2$ 时, 由上式得:
$$R(n) \leq 2c + \frac{1}{2} \max\{R(1), R(1)\} \\ \leq 2.5c < 4cn$$

② 假设对所有的 $n, 2 \leq n < m$, 有 $R(n) \leq 4cn$

③ 当 $n = m$ 时, 有,
$$R(n) \leq cm + \frac{1}{m} \max_k \left\{ \sum_{i=m-k+1}^{m-1} R(i) + \sum_{i=k}^{m-1} R(i) \right\}$$

由于 $R(n)$ 是 n 的非降函数, 故在当 m 为偶数而 $k=m/2$, 或当 m 为奇数而 $k=(m+1)/2$ 时, $\sum_{i=n-k+1}^{n-1} R(i) + \sum_{i=k}^{n-1} R(i)$ 取得极大值。因此,

若 m 为偶数, 则
$$R(m) \leq cm + \frac{2}{m} \sum_{i=m/2}^{m-1} R(i) \leq cm + \frac{8c}{m} \sum_{i=m/2}^{m-1} i < 4cm$$

若 m 为奇数, 则
$$R(m) \leq cm + \frac{2}{m} \sum_{i=(m+1)/2}^{m-1} R(i) \leq cm + \frac{8c}{m} \sum_{i=(m+1)/2}^{m-1} i < 4cm$$

由于 $T_A(n) \leq R(n)$, 所以 $T_A(n) \leq 4cn$ 。故, $T_A(n) = O(n)$

最坏情况是 $O(n)$ 的选择算法

1) 采用两次取中间值的规则精心选取划分元素

首先, 将参加划分的 n 个元素分成 $\lfloor n/r \rfloor$ 组, 每组有 r 个元素($r \geq 1$)。 (多余的 $n - r\lfloor n/r \rfloor$ 个元素忽略不计)

然后, 对这 $\lfloor n/r \rfloor$ 组每组的 r 个元素进行分类并找出其中间元素 $m_i, 1 \leq i \leq \lfloor n/r \rfloor$, 共得 $\lfloor n/r \rfloor$ 个中间值。

之后, 对这 $\lfloor n/r \rfloor$ 个中间值分类, 并找出其中间值 m_m 。
将 m_m 作为划分元素执行划分。

2) 算法描述

算法 使用二次取中规则选择算法 SELECT2(A,k,n)

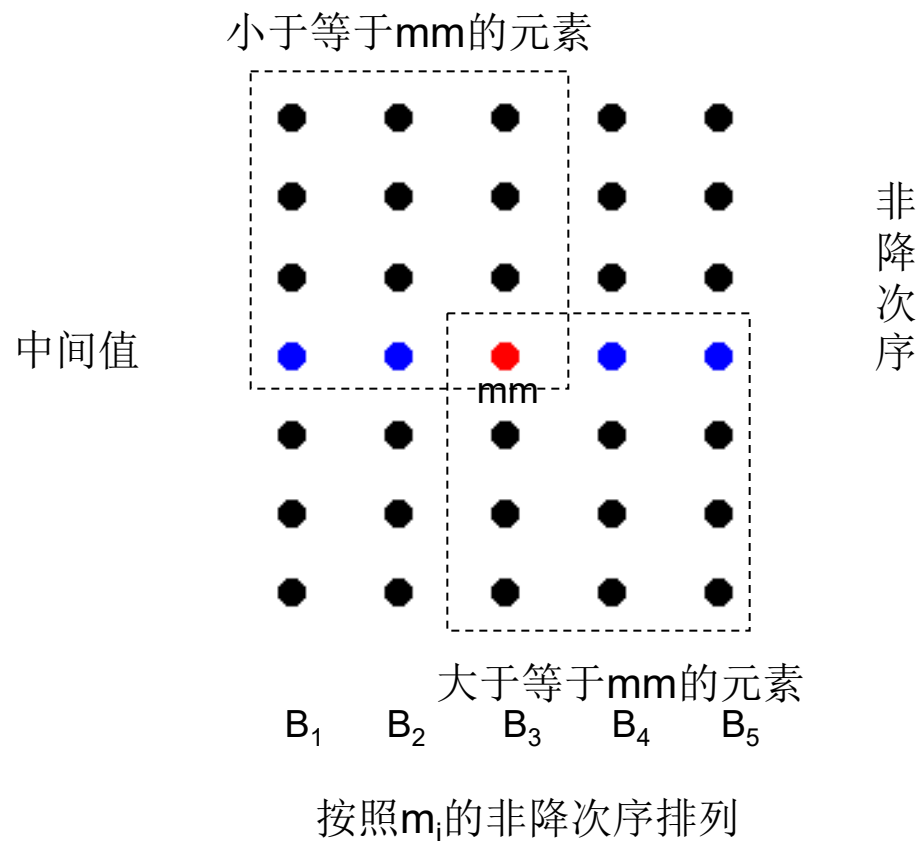
//在集合A中找第k小元素，使用两次取中规则//

- ① 若 $n \leq r$ ，则采用插入法直接对A分类并返回第k小元素
 - ② 把A分成大小为r的 $\lfloor n/r \rfloor$ 个子集合，忽略多余的元素
 - ③ 设 $M = \{m_1, m_2, \dots, m_{\lfloor n/r \rfloor}\}$ 是 $\lfloor n/r \rfloor$ 子集合的中间值集合
 - ④ $v \leftarrow \text{SELECT2}(M, \lceil \lfloor n/r \rfloor / 2 \rceil, \lfloor n/r \rfloor)$
 - ⑤ $j \leftarrow \text{PARTITION}(A, v)$ //v作为划分元素，划分后j等于划分元素所在位置的下标//
 - ⑥ case
 - : $k=j$: return(v)
 - : $k < j$: 设S是A(1:j-1)中元素的集合
return(SELECT2(S,k,j-1))
 - :else: 设R是A(j+1:n)中元素的集合
return(SELECT2(R,k-j,n-j))endcase
- end SELECT2

• 例：设 $n=35$, $r=7$ 。

- 分为 $n/r = 5$ 个元素组：B1, B2, B3, B4, B5；每组有7个元素。
- B1-B5按照各组的 m_i 的非增次序排列。
- $mm = m_i$ 的中间值, $1 \leq i \leq 5$

由图所示有：



由于 r 个元素的中间值是第 $\lceil r/2 \rceil$ 小元素。则,

- 至少有 $\lceil \lfloor n/r \rfloor / 2 \rceil$ 个 m_i 小于或等于 m_m ;
- 至少有 $\lfloor n/r \rfloor - \lceil \lfloor n/r \rfloor / 2 \rceil + 1 \geq \lceil \lfloor n/r \rfloor / 2 \rceil$ 个 m_i 大于或等于 m_m 。

则, 至少有 $\lceil r/2 \rceil \lceil \lfloor n/r \rfloor / 2 \rceil$ 个元素小于或等于 (或大于或等于) m_m 。

- 当 $r=5$, 则使用两次取中间值规则来选择 $v=m_m$, 可推出,
- 至少有 $1.5 \lfloor n/5 \rfloor$ 个元素小于或等于选择元素 v 。
- 至多有 $n - 1.5 \lfloor n/5 \rfloor \leq 0.7n + 1.2$ 个元素大于 v 。
- 至多有 $0.7n + 1.2$ 个元素小于 v 。

故, 这样的 v 可近似平均地划分 A 中的 n 个元素。

记 $T(n)$ 是SELECT2所需的最坏情况时间

对特定的 r 分析SELECT2:选取 $r=5$ 。

➤ 假定 A 中的元素各不相同，则有

① 若 $n \leq r$ ，则采用插入法直接对 A 分类并返回第 k 小元素 $\rightarrow O(1)$

② 把 A 分成大小为 r 的 $\lfloor n/r \rfloor$ 子集合，忽略多余的元素 $\rightarrow O(n)$

③ 设 $M = \{m_1, m_2, \dots, m_{\lfloor n/r \rfloor}\}$ 是 $\lfloor n/r \rfloor$ 子集合的中间值集合 $\rightarrow O(n)$

④ $v \leftarrow \text{SELECT2}(M, \lfloor \lfloor n/r \rfloor / 2 \rfloor)$ $\rightarrow T(n/5)$

⑤ $j \leftarrow \text{PARTITION}(A, v)$ $\rightarrow O(n)$

⑥ case $\rightarrow T(3n/4), n \geq 24$

: $k=j$: return(v)

: $k < j$: 设 S 是 $A(1:j-1)$ 中元素的集合; return(SELECT2($S, k, j-1$))

:else: 设 R 是 $A(j+1:n)$ 中元素的集合; return(SELECT2($R, k-j, n-j$))

endcase

故有,

$$T(n) = \begin{cases} cn & n < 24, \\ T(n/5) + T(3n/4) + cn & n \geq 24 \end{cases}$$

用归纳法可证:

$$T(n) \leq 20cn$$

$$\text{即, } T(n) = O(n)$$

➤ 当A中的元素不尽相同

步骤⑤经PARTITION调用所产生的S和R两个子集合中可能存在一些元素等于当前的划分元素 v ,可能导致 $|S|$ 或 $|R|$ 大于 $0.7n+1.2$ 。

此时上述处理 作一下改进:

方法一: 将A集合分成3个子集合U,S和R, 其中U是有A中所有与 v 相同的元素组成, S是由A中所有比 v 小的元素组成, R则是A中所有比 v 大的元素组成。

同时步骤⑥更改:

case

: $|S| \geq k$: return(SELECT2(S,k, $|S|$))

: $|S| + |U| \geq k$: return(v)

:else: return(SELECT2(R,k- $|S|$ - $|U|$, $|R|$))

endcase

$$T(n) = O(n)$$

从而保证 $|S| + |R| \leq 0.7n + 1.2$ 成立, 故关于 $T(n)$ 的分析仍然成立。

方法二：选取其他的 r 值进行计算

特例：当 $r=5$,且 A 中的元素不尽相同。假设其中有 $0.7n+1.2$ 个元素比 v 小而其余的元素都等于 v 的情况。

则，经过PARTITION，在这些等于 v 的元素中至多有一半可能在 S 中，故
 $|S| \leq 0.7n+1.2+(0.3n-1.2)/2=0.85n+0.6$

同理， $|R| \leq 0.85n+0.6$

可得，步骤④和⑥此时所处理的元素总数将是

$$1.05n+0.6 > n$$

不再是线性关系。故有 $T(n) \neq O(n)$

改进：

取 $r=9$ 。经计算可得，此时将有 $2.5 \lfloor n/9 \rfloor$ 个元素小于或等于 v ，同时至少有同样多的元素大于或等于 v 。

则当 $n \geq 90$ 时， $|S|$ 和 $|R|$ 都至多为

$$n - 2.5 \lfloor n/9 \rfloor + \frac{1}{2} (2.5 \lfloor n/9 \rfloor) = n - 1.25 \lfloor n/9 \rfloor \leq 31n/36 + 1.25 \leq 63n/72$$

故有,

$$T(n) = \begin{cases} c_1 n & n < 90 \\ T(n/9) + T(63n/72) + c_1 n & n \geq 90 \end{cases}$$

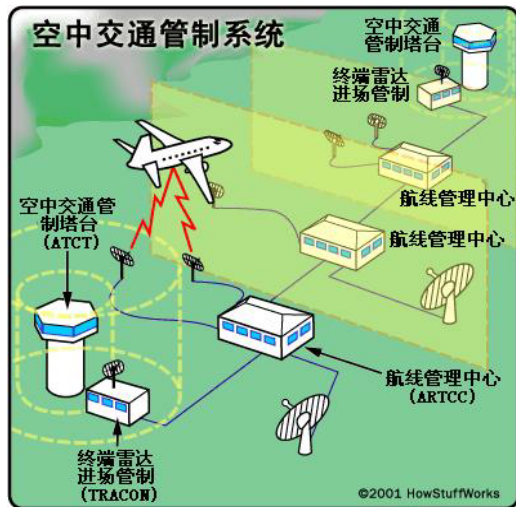
用归纳法可证:

$$T(n) \leq 72c_1 n$$

即, $T(n) = O(n)$

分治-例7 最接近点对问题

- 给定平面上 n 个点的集合 S ，找其中的一对点，使得在 n 个点组成的所有点对中，该点对间距离最小。
 - 应用：空中交通管制



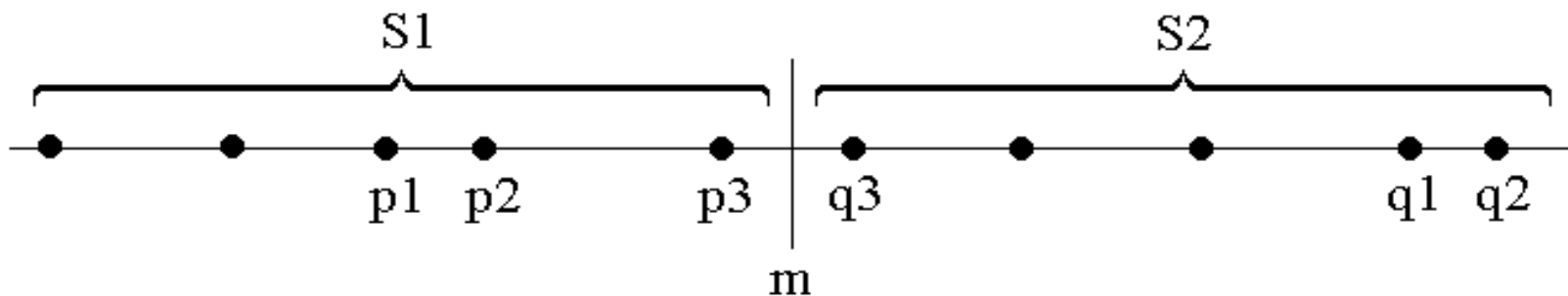
直接求解方法

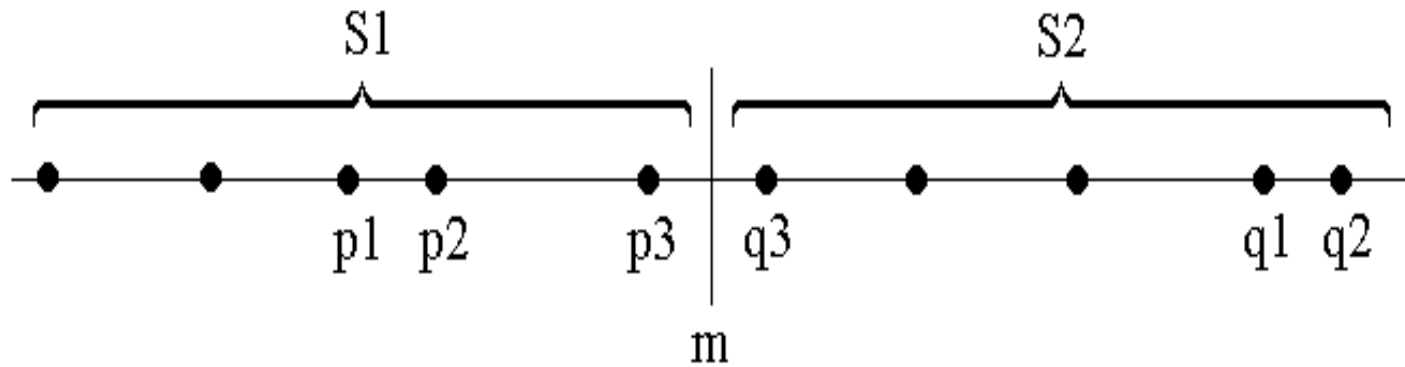
- 分别计算每一对点之间的距离，然后找出距离最小的那一对。
- 时间复杂度

$$\begin{aligned}C(n) &= \sum_{i=1}^{n-1} \sum_{j=i+1}^n 2 = 2 \sum_{i=1}^{n-1} n - i \\&= 2[(n-1) + (n-2) + \dots + 1] \\&= (n-1)n \in \Theta(n^2)\end{aligned}$$

◆为了使问题易于理解和分析，先来考虑**一维**的情形。
此时， S 中的 n 个点退化为 x 轴上的 n 个实数 x_1, x_2, \dots, x_n 。
最接近点对即为这 n 个实数中相差最小的2个实数。

- 假设我们用 x 轴上某个点 m 将 S 划分为2个子集 S_1 和 S_2 ，基于**平衡子问题**的思想，用 S 中各点坐标的中位数来作分割点。
- 递归地在 S_1 和 S_2 上找出其最接近点对 $\{p_1, p_2\}$ 和 $\{q_1, q_2\}$ ，并设 $d = \min\{|p_1 - p_2|, |q_1 - q_2|\}$ ， S 中的最接近点对或者是 $\{p_1, p_2\}$ ，或者是 $\{q_1, q_2\}$ ，或者是某个 $\{p_3, q_3\}$ ，其中 $p_3 \in S_1$ 且 $q_3 \in S_2$ 。
- 能否在线性时间内找到 p_3, q_3 ？





能否在线性时间内找到 p_3, q_3 ?

- ◆ 如果 S 的最接近点对是 $\{p_3, q_3\}$, 即 $|p_3 - q_3| < d$, 则 p_3 和 q_3 两者与 m 的距离不超过 d , 即 $p_3 \in (m-d, m]$, $q_3 \in (m, m+d]$.
- ◆ 由于在 S_1 中, 每个长度为 d 的半闭区间至多包含一个点 (否则必有两点距离小于 d), 并且 m 是 S_1 和 S_2 的分割点, 因此 $(m-d, m]$ 中至多包含 S 中的一个点。由图可以看出, 如果 $(m-d, m]$ 中有 S 中的点, 则此点就是 S_1 中最大点。
- ◆ 因此, 我们用线性时间就能找到区间 $(m-d, m]$ 和 $(m, m+d]$ 中所有点, 即 p_3 和 q_3 。从而我们用线性时间就可以将 S_1 的解和 S_2 的解合并成为 S 的解。

```
public static double cpair1(S)
```

```
{    n=|S|;
```

```
    if (n < 2) return  $\infty$ ;
```

```
    m=S中各点坐标的中位数; ;
```

```
    d1=cpair1(S1); // S1={x  $\in$  S|x $\leq$ m}, 构造S1和S2
```

```
    d2=cpair1(S2); // S2={x  $\in$  S|x>m}, 构造S2
```

```
    p=max(S1);
```

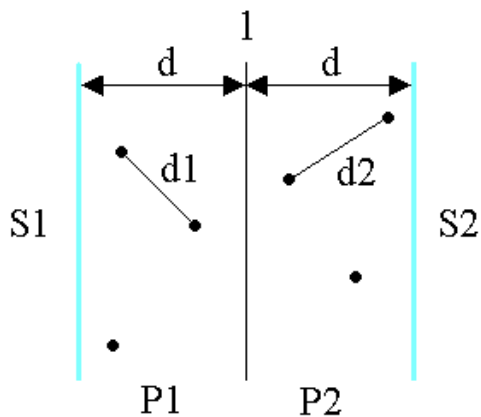
```
    q=min(S2);
```

```
    d=min(d1,d2,q-p);
```

$$T(n) = \begin{cases} O(1) & n < 4 \\ 2T(n/2) + O(n) & n \geq 4 \end{cases}$$

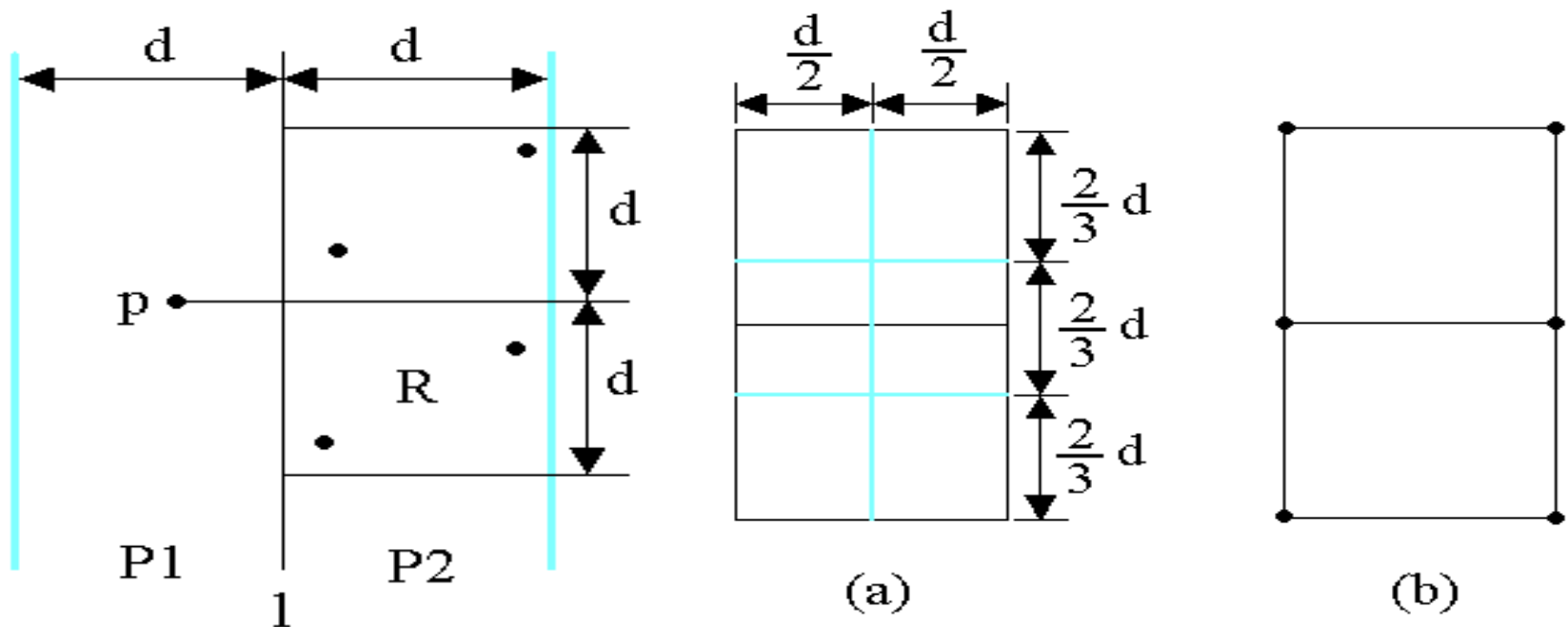
二维情形

- ◆ 下面来考虑二维的情形。
- ◆ 选取一垂直线 $l: x = m$ 来作为分割直线。其中 m 为 S 中各点 x 坐标的中位数。由此将 S 分割为 S_1 和 S_2 。
- ◆ 递归地在 S_1 和 S_2 上找出其最小距离 d_1 和 d_2 ，并设 $d = \min\{d_1, d_2\}$ ， S 中的最接近点对或者是 d ，或者是某个 $\{p, q\}$ ，其中 $p \in P_1$ 且 $q \in P_2$ 。
- ◆ 能否在线性时间内找到 p, q ？



二维情形

- 考虑 P_1 中任意一点 p ，它若与 P_2 中的点 q 构成最接近点对的候选者，则必有 $\text{distance}(p, q) < d$ 。满足这个条件的 P_2 中的点一定落在一个 $d \times 2d$ 的矩形 R 中
- 由 d 的意义可知， P_2 中任何2个 S 中的点的距离都不小于 d 。由此可以推出矩形 R 中最多只有6个 S 中的点。
- 因此，在分治法的合并步骤中最多只需要计算 $n/2 \times 6$ 个点对。



二维情形



证明:将矩形R的长为 $2d$ 的边3等分, 将它的长为 d 的边2等分, 由此导出6个 $(d/2) \times (2d/3)$ 的矩形。若矩形R中有多于6个S中的点, 则由鸽舍原理易知至少有一个 $(d/2) \times (2d/3)$ 的小矩形中有2个以上S中的点。设 u, v 是位于同一小矩形中的2个点, 则:

$$\begin{aligned}(x(u) - x(v))^2 + (y(u) - y(v))^2 &\leq (d/2)^2 + (2d/3)^2 \\ &\leq \frac{25}{36}d^2\end{aligned}$$

$\text{distance}(u,v) < d$ 。这与 d 的意义相矛盾。

二维情形

➤ 要检查哪6个点？

将 p 和 P_2 中所有 S_2 的点投影到垂直线 l 上。由于能与 p 点一起构成最接近点对候选者的 S_2 中点一定在矩形 R 中，所以它们在直线 l 上的投影点距 p 在 l 上投影点的距离小于 d ，这种投影点最多只有6个。

➤ 因此，若将 P_1 和 P_2 中所有 S 中点按其 y 坐标排好序，则对 P_1 中所有点，对排好序的点列作一次扫描，就可以找出所有最接近点对的候选者。对 P_1 中每一点最多只要检查 P_2 中排好序的相继6个点。

double **cpair2**(S)

{

n=|S|;

if (n < 2) **return** ;

1、 m=S中各点x|间坐标的
中位数;构造S1和S2;

$S1 = \{p \in S | x(p) \leq m\}$,

$S2 = \{p \in S | x(p) > m\}$

2、 d1=**cpair2**(S1);

d2=**cpair2**(S2);

3、 dm=**min**(d1,d2);

4、 设P1 是S1中距垂直分割线 l 的距离在
dm 之内的所有点组成的集合;

P2 是S2 中距分割线 l 的距离在dm 之
内所有点组成的集合;

将P1 和P2 中点依其 y 坐标值排序;

并设X 和Y 是相应的已排好序的点列;

5、 通过扫描 X 以及对于 X 中每个点检查
Y中与其距离在dm 之内的所有点(最多6个)
可以完成合并;

当X中的扫描指针逐次向上移动时, Y
中的扫描指针可在宽为2dm的区间内移动;

设dl 是按这种扫描方式找到的点对间
的最小距离;

6、 d=**min**(dm,dl);

return d;

}

2.10 最接近点对问题

$$T(n) = \begin{cases} O(1) & n < 4 \\ 2T(n/2) + O(n) & n \geq 4 \end{cases}$$

复杂度分析

$n \geq 4$,

$$\begin{aligned} T(n) &= 2T(n/2) + O(n) = 2[2T(n/4) + O(n/2)] + O(n) \\ &= 2^2 T(n/2^2) + 2O(n) = \dots \end{aligned}$$

$$n = 2^k \quad = 2^k + kO(n)$$

$$k \geq 2 \quad \mathbf{T(n) = O(n \log n)}$$

分治-例8 循环赛日程表

- 有 $n=2k$ 个运动员。设计一个满足以下要求的比赛日程表：
 - (1)每个选手必须与其他 $n-1$ 个选手各赛一次；
 - (2)每个选手一天只能赛一次；
 - (3)循环赛一共进行 $n-1$ 天。

思想

- 按分治策略，将所有的选手分为两半
 - n 个选手的比赛日程表就可以通过为 $n/2$ 个选手设计的比赛日程表来决定。
- 递归地对选手进行分割，直到只剩下2个选手时，只要让这2个选手进行比赛就可以了。

示例

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