



# 随机算法

朱同鑫

# 本章内容



- 随机算法的基本概念
- 随机数值算法
  - 计算 $\pi$ 值、计算定积分
- Las Vegas算法
  - 第k小元素问题
- Sherwood随机化方法
- Monte Carlo算法
  - 字符串相等判断问题
- 子串匹配问题
- 最近点对随机算法
- 素数测试



# 随机算法的基本概念

- 从一个例子看随机算法的作用

- 【例】判断函数  $f(x_1, x_2, \dots, x_n)$  在区域  $D$  中是否恒为0,  $f$  很复杂, 不能数学简化, 给出确定性结果很麻烦
  - 若随机产生一个  $n$  维坐标  $(r_1, r_2, \dots, r_n) \in D$ , 代入得  $f(r_1, r_2, \dots, r_n) \neq 0$ , 则可判定区域  $D$  内  $f$  不恒为0
  - 若对**很多个**随机产生的坐标进行测试, 结果次次均为0, 则可说  $f \neq 0$  的概率是**非常小** (测试数量越多, 可信度越高)



**结论:** 有不少问题, 目前只有效率很差的确定性求解算法, 但用随机算法去求解, 可以很快地获得相当可信的结果。



# 随机算法的基本概念

- 随机算法

- 随机算法是一种使用概率和统计方法在其执行过程中对于下一计算步骤作出随机选择的算法

- 随机算法的优越性

- 对于有些问题：算法简单
- 对于有些问题：时间复杂性低
- 对于有些问题：同时兼有简单和时间复杂性低



# 随机算法的基本概念

- 随机算法的随机性
  - 对于同一实例的多次执行, 效果可能完全不同
  - 时间复杂性是一个随机变量
  - 解的正确性和准确性也是随机的
- 随机算法的分类
  - 随机数值算法
  - Monte Carlo算法
  - Las Vegas算法
  - Sherwood算法



# 随机算法的基本概念

- 随机数值算法
  - 主要用于数值问题求解
  - 算法的输出往往是近似解
  - 近似解的精确度与算法执行时间成正比
- Monte Carlo算法
  - 主要用于求解需要准确解的问题
  - 算法可能给出错误的解
  - 获得正确解的概率与算法执行时间成正比



# 随机算法的基本概念

- Las Vegas算法
  - 可能找不到解，但一旦找到一个解，该解一定是正确的
  - 找到解的概率与算法执行时间成正比
  - 增加对问题反复求解次数, 可使求解无效的概率任意小
- Sherwood算法
  - 一定能够求得一个正确解
  - 确定算法的最坏与平均复杂性差别大时，加入随机性，即得到Sherwood算法
  - 消除最坏行为与特定实例的联系



# 随机算法的基本概念

- 随机算法分析的特征
  - 仅依赖于随机选择：不依赖于输入的分布
  - 确定算法的平均复杂性分析：依赖于输入的分布
- 随机算法分析的目标
  - 平均时间复杂性：时间复杂性随机变量的期望
  - 获得正确解的概率/获得优化解的概率
  - 解的精确度分析

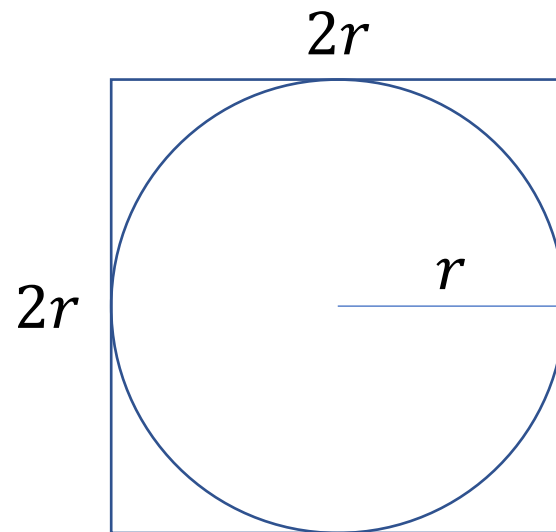




# 随机数值算法——计算 $\pi$ 值

- 数学基础

- 设有一个半径为  $r$  的圆及其外切四边形
- 向正方形随机地投掷点，投掷点落入圆内的概率为  $\frac{\pi r^2}{4r^2} = \frac{\pi}{4}$
- 投掷  $n$  个点，设  $k$  个点落入圆内
- 用  $k/n$  逼近  $\pi/4$ ，即  $k/n \approx \pi/4$ ，于是  $\pi \approx 4k/n$





# 随机数值算法——计算 $\pi$ 值

- 随机算法伪代码

## Computing $\pi(n)$

$k = 0;$

**for**  $i = 1$  **to**  $n$  **do**

    随机地产生正方形中的一点 $(x, y);$

**if**  $x^2 + y^2 \leq 1$  **then**

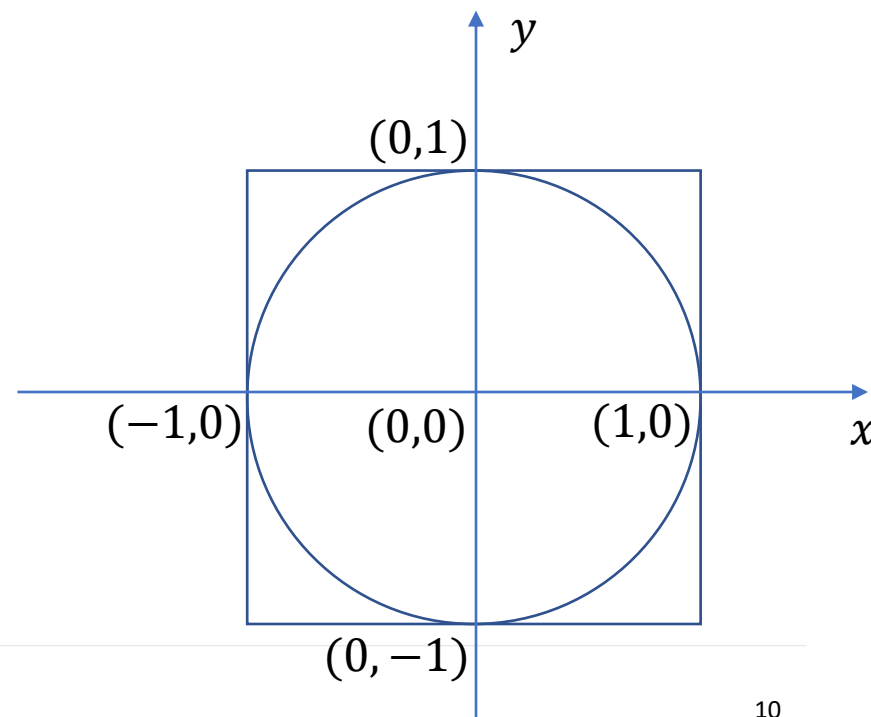
$k = k + 1;$

**return**  $4k/n;$

解的精确度随着随机样本大小 $n$   
增加而增加

 **时间复杂度 $O(n)$**

$n$ 不是输入的大小，而是随机样本的大小





# 随机数值算法——计算定积分

- 问题：计算积分  $\int_a^b g(x)dx$
- 数学基础
  - **强大数定律**：假定  $\{s(x)\}$  是相互独立同分布的随机变量序列，如果它们有有限的数学期望  $E(s(x)) = a$ ，则

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n s(x_i) = a = E(s(x))$$



# 随机数值算法——计算定积分

## • 计算积分

- 令  $f(x) = \frac{1}{b-a}$  为区间  $[a, b]$  上的一组独立、同分布的随机变量  $\{s(x)\}$  的概率密度函数
  - 当  $\{s(x)\}$  为离散随机变量时, 期望  $E(s(x)) = \sum s(x)f(x)$
  - 当  $\{s(x)\}$  为连续随机变量时, 期望  $E(s(x)) = \int_a^b s(x)f(x) dx$
- 令  $g(x) = s(x)f(x)$ , 那么期望  $E(s(x)) = \int_a^b g(x) dx$
- 则,  $\int_a^b g(x) dx = E(s(x)) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n s(x_i) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n g(x_i)/f(x_i)$   
$$= \lim_{n \rightarrow \infty} \frac{b-a}{n} \sum_{i=1}^n g(x_i)$$



# 随机数值算法——计算定积分

- 随机算法伪代码

## Computing-Definite-Integration $(a, b, n)$

$R = 0;$

**for**  $i = 1$  **to**  $n$  **do**

    随机地产生  $[a, b]$  中的点  $X$ ;

$R = R + g(x);$

**return**  $(b - a)R/n;$

$$\int_a^b g(x) dx = \lim_{n \rightarrow \infty} \frac{b-a}{n} \sum_{i=1}^n g(x_i)$$

解的精确度随着随机样本大小  $n$  增加而增加

 **时间复杂度  $O(n)$**

$n$  不是输入的大小，而是随机样本的大小



# Las Vegas算法——第k小元素问题

- 问题定义

- 输入:  $S = \{x_1, x_2, \dots, x_n\}$ , 整数  $k$ ,  $1 \leq k \leq n$
- 输出:  $S$  中第  $k$  小元素

- 随机算法

- 在  $n$  个数中随机找一个数  $x_i$ , 然后将剩余  $n - 1$  个数与  $x_i$  比价, 分别放入三个数组  $S_1$  (元素均  $< x_i$ )、 $S_2$  (元素均  $= x_i$ )、 $S_3$  (元素均  $> x_i$ ) 中
  - 当  $k \leq |S_1|$  时, 调用  $Select(S_1, k)$
  - 当  $|S_1| < k \leq |S_1| + |S_2|$  时, 返回  $x_i$
  - 当  $k > |S_1| + |S_2|$  时, 调用  $Select(S_3, k - |S_1| - |S_2|)$



# Las Vegas算法——第k小元素问题

## • 性能分析

### 定理.

若以等概率方法在 $n$ 个数中随机取数，则该算法用到的比较次数的期望值不超过 $4n$ 。

### 证明.

设 $C(n)$ 是输入规模为 $n$ 时，算法比较次数的期望值，并设取到任意数的概率相同，假设取到第 $j$ 小的数。

- 若 $j > k$ ，需要调用 $Select(S_1, k)$ ，而 $|S_1| = j - 1$ ，比较次数期望为 $C(j - 1)$
- 若 $j = k$ ，无需继续进行比较，直接返回第 $j = k$ 个元素
- 若 $j < k$ ，需要调用 $Select(S_3, k - |S_1| - |S_2|)$ ，而 $|S_3| = n - j$ ，比较次数期望为 $C(n - j)$



# Las Vegas算法——第k小元素问题

## • 性能分析

### 定理.

若以等概率方法在 $n$ 个数中随机取数，则该算法用到的比较次数的期望值不超过 $4n$ 。

### 证明.

$$\begin{aligned} C(n) &= n + \frac{1}{n} \left( \sum_{j=k+1}^n C(j-1) + \sum_{j=1}^{k-1} C(n-j) \right) \\ &= n + \frac{1}{n} \left( \sum_{i=k}^{n-1} C(i) + (C(n-1) + C(n-2) + \cdots + C(n-k+1)) \right) \\ &= n + \frac{1}{n} \left( \sum_{i=k}^{n-1} C(i) + \sum_{i=n-k+1}^{n-1} C(i) \right) \\ &\leq n + \frac{1}{n} \left( \sum_{i=n-\frac{n}{2}+1}^{n-1} C(i) + \sum_{i=\frac{n}{2}}^{n-1} C(i) \right) \end{aligned}$$

由于 $C(i)$ 是非减函数，即 $i < j$ 时，总有 $C(i) \leq C(j)$ ， $C(n)$ 在 $k = \lceil n/2 \rceil$ 时取得最大值





# Las Vegas算法——第k小元素问题

## • 性能分析

### 定理.

若以等概率方法在 $n$ 个数中随机取数，则该算法用到的比较次数的期望值不超过 $4n$ 。

### 证明.

归纳法证明 $C(n) \leq n + \frac{1}{n} \left( \sum_{i=n-\frac{n}{2}+1}^{n-1} C(i) + \sum_{i=\frac{n}{2}}^{n-1} C(i) \right) \leq 4n$

1. 当 $n = 1$ 时,  $C(1) \leq 4$ 显然成立
2. 假设当 $n \leq m - 1$ 时 ( $m \geq 2$ ) ,  $C(n) \leq 4n$ 成立, 往证当 $n = m$ 时,  $C(n) \leq 4n$ 成立

$$C(m) \leq m + \frac{1}{m} \left( \sum_{i=m-\frac{m-1}{2}+1}^{m-1} 4i + \sum_{i=\frac{m}{2}}^{m-1} 4i \right) \leq m + \frac{1}{m} \left( 8 \sum_{i=\frac{m}{2}}^{m-1} i \right) \leq 4m$$



# Sherwood随机化方法

- 一般过程

- 若问题已经有平均性质较好的确定性算法，但是该算法在最坏情况下效率不高
- 引入一个随机数发生器（通常服从均匀分布），将一个确定性算法改成一快个随机算法
- 【例】速排序
  - 每次选择一个基准数，比它小的放左边，大的放右边，如此递归
  - 平均效率很好，但是最坏情况下，每次选择的基准若都是最小的，或是最大的，算法效率不高
  - 可随机预处理(洗牌)，使输入均匀分布，再运行算法



# Monte Carlo算法——字符串相等判断

- 问题描述:

- 设A处有一个长字符串 $X$  (e.g. 长度为 $10^6$ ) , B处也有一个长字符串 $Y$  , A将 $X$ 发给B, 由B判断是否有 $X = Y$

- 判断方法:

- 首先, A将 $X$ 的长度发给B, 若B判断该长度与 $Y$ 的长度不等, 则 $X \neq Y$
- 若长度相等, 采用 “取指纹” 的判断方法
  - A对 $X$ 进行处理, 取出 $X$ 的 “指纹” , 将指纹发给B
  - 由B检查 $X$ 的 “指纹” 是否等于 $Y$ 的 “指纹”
  - 若取 $k$ 次指纹(每次指纹不同), 每次两者结果均相同, 则认为 $X = Y$ 。
  - 随着 $k$ 的增大, 误判率可趋于0



# Monte Carlo算法——字符串相等判断

- 字符串取“指纹”方法

- 令字符串 $X$ 的二进制编码长度为 $n$ , 则  $X < 2^n$ , 取  $f(X) \equiv X \pmod{p}$  作为  $X$  的指纹, 其中  $p$  是小于  $2n^2$  的素数,
- $p$  的二进制长度:  $\log_2 p \leq \lfloor \log_2 2n^2 \rfloor + 1 = O(\log_2 n)$ ,  $f(X)$  的二进制长度  $\leq \log_2 p$ , 即传输长度可大大缩短
- 【例】 $X$  的二进制是  $10^6$  位, 即  $n = 10^6$ , 则  $p < 2 \times 10^{12} \approx 2^{40.8631}$ 
  - $f(X)$  的位数不超过 41 位, 传输一次指纹  $f(X)$  可节省约 2.5 万倍
  - 根据下面所做的分析, 错判率小于  $\frac{1}{10^6}$
  - 若取 5 次指纹, 错判率小于  $\frac{1}{10^{30}}$



# Monte Carlo算法——字符串相等判断

- 错判率分析

- B接到指纹 $f(X)$ 后与 $f(Y)$ 比较,
  - 若 $f(X) \neq f(Y)$ , 当然有 $X \neq Y$ 。
  - 若 $f(X) = f(Y)$ 而 $X \neq Y$ , 此时是一个错误匹配。
- 错误匹配的概率有多大?
  - 随机取一个小于 $2n^2$ 的素数 $p$ , 则对于给定的 $X$ 和 $Y$ ,
  - 错判率

$$P_{fail} = \frac{X \neq Y, \text{但使得 } f(X) = f(Y) \text{ 的小于 } 2n^2 \text{ 的素数的个数}}{\text{小于 } 2n^2 \text{ 的素数的总个数}}$$



# Monte Carlo算法——字符串相等判断

## • 错判率分析

• 错判率:  $P_{fail} = \frac{X \neq Y, \text{但使得 } f(X) = f(Y) \text{ 的小于 } 2n^2 \text{ 的素数的个数}}{\text{小于 } 2n^2 \text{ 的素数的总个数}} \leq \frac{\pi(n)}{\pi(2n^2)} \leq \frac{1}{n}$

数论定理1: 设 $\pi(n)$ 是小于 $n$ 的素数个数, 则 $\pi(n) \approx \frac{n}{\ln n}$

则小于 $2n^2$ 的素数的总个数为:  $\pi(2n^2) \approx \frac{n^2}{\ln n}$

数论定理2: 设 $a \equiv b \pmod{p}$  iff  $p$ 能整除 $|a - b|$

使得 $f(X) = f(Y)$ 的素数的个数=能够整除 $|X - Y|$ 素数的个数

数论定理3: 若 $a < 2^n$ , 则能整除 $a$ 的素数个数不超过 $\pi(n)$ 个

$|X - Y| < \max\{X, Y\} \leq 2^n - 1$ , 所以能整除 $|X - Y|$ 的素数个数  $\leq \pi(n)$



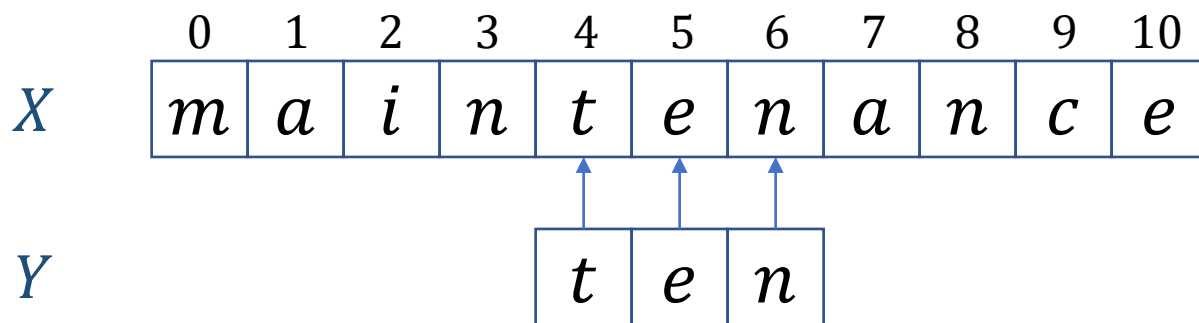
# 子串匹配问题

- 问题定义:

- 输入: 目标串  $X = "x_1x_2 \cdots x_n"$ , 模式串  $Y = "y_1y_2 \cdots y_m"$ ,  $m \leq n$
- 输出: 判断模式串  $Y$  是否为目标串  $X$  的子串

- 问题实例

- $X = "maintenance"$ ,  $Y = "ten"$



**$Y$  是  $X$  的子串, 且在  $X$  中的位置为 4**



# 子串匹配问题

- 随机算法

- 记  $X(j) = x_j x_{j+1} \cdots x_{j+m-1}$
- 逐一比较  $X(j)$  的指纹  $f(X(j)) \equiv X(j) \pmod{p}$  与  $Y$  的指纹  $f(Y) \equiv Y \pmod{p}$ ,  
 $1 \leq j \leq n - m + 1$ ,  $p$  是小于  $2mn^2$  的素数
- $f(X(j+1))$  可根据  $f(X(j))$  计算, 故算法可很快完成





# 子串匹配问题

## • 随机算法

数论定理4:  $(xy + z) \pmod p = (x(y \pmod p) + z) \pmod p$

- $f(X(j+1))$ 可根据 $f(X(j))$ 计算
  - 不失一般性, 令 $X$ 和 $Y$ 都是0-1串 (二进制编码)
  - $f(X(j+1)) = (x_{j+1} \cdots x_{j+m-1}x_{j+m}) \pmod p$
  - $= (2(x_{j+1} \cdots x_{j+m-1}) + x_{j+m}) \pmod p$
  - $= (2(x_{j+1} \cdots x_{j+m-1}) + 2^m x_j - 2^m x_j + x_{j+m}) \pmod p$
  - $= (2(x_j x_{j+1} \cdots x_{j+m-1}) - 2^m x_j + x_{j+m}) \pmod p$
  - $= (2((x_j x_{j+1} \cdots x_{j+m-1}) \pmod p) - 2^m x_j + x_{j+m}) \pmod p$
  - $= (2f(X(j)) - (2^m \pmod p)x_j + x_{j+m}) \pmod p$

被余数都在 $p$ 附近, 与 $p$ 相差很小 ( $[-(p-1), 2p-1]$ ), 只需一两次加减法即可求余



# 子串匹配问题

## • 错判率分析

- 当 $Y \neq X(j)$ , 但 $f(Y) = f(X(j))$ 时产生错误
- 而 $f(Y) = f(X(j))$ 当且仅当 $p$ 能整除 $|Y - X(j)|$ 
  - 若 $p$ 能整除 $|Y - X(j)|$ ,  $p$ 自然也能整除 $Z = |Y - X(1)| \times \cdots \times |Y - X(n - m + 1)|$
  - 因为 $|Y - X(j)| \leq 2^m$ ,
  - 所以 $Z = |Y - X(1)| \times \cdots \times |Y - X(n - m + 1)| \leq (2^m)^{n-m+1} \leq 2^{mn}$

数论定理3: 若 $a < 2^n$ , 则能整除 $a$ 的素数个数不超过 $\pi(n)$ 个

- 所以能整除 $Z$ 的素数个数不超过 $\pi(mn)$ 个

$$P_{fail} = \frac{X \text{ 不包含 } Y, \text{ 但能整除 } Z \text{ 的素数的个数}}{\text{小于 } 2^{mn^2} \text{ 的素数的总个数}} \leq \frac{\pi(mn)}{\pi(2^{mn^2})} = \frac{1}{n}$$



# 子串匹配问题

- 确定算法1——穷举模式匹配算法
  - 匹配失败，目标串 $X$ 回溯，模式串 $Y$ 从头开始

	$X$	<table><tr><td><math>x_1</math></td><td><math>x_2</math></td><td><math>x_3</math></td><td><math>x_4</math></td><td><math>\cdots</math></td><td><math>x_{m-2}</math></td><td><math>x_{m-1}</math></td><td><math>x_m</math></td><td><math>\cdots</math></td><td><math>x_n</math></td></tr></table>	$x_1$	$x_2$	$x_3$	$x_4$	$\cdots$	$x_{m-2}$	$x_{m-1}$	$x_m$	$\cdots$	$x_n$
$x_1$	$x_2$	$x_3$	$x_4$	$\cdots$	$x_{m-2}$	$x_{m-1}$	$x_m$	$\cdots$	$x_n$			
第1趟	$Y$	<table><tr><td><math>y_1</math></td><td><math>y_2</math></td><td><math>y_3</math></td><td><math>y_4</math></td><td><math>\cdots</math></td><td><math>y_{m-2}</math></td><td><math>y_{m-1}</math></td><td><math>y_m</math></td></tr></table>	$y_1$	$y_2$	$y_3$	$y_4$	$\cdots$	$y_{m-2}$	$y_{m-1}$	$y_m$		
$y_1$	$y_2$	$y_3$	$y_4$	$\cdots$	$y_{m-2}$	$y_{m-1}$	$y_m$					
第2趟	$Y$	<table><tr><td><math>y_1</math></td><td><math>y_2</math></td><td><math>y_3</math></td><td><math>y_4</math></td><td><math>\cdots</math></td><td><math>y_{m-2}</math></td><td><math>y_{m-1}</math></td><td><math>y_m</math></td></tr></table>	$y_1$	$y_2$	$y_3$	$y_4$	$\cdots$	$y_{m-2}$	$y_{m-1}$	$y_m$		
$y_1$	$y_2$	$y_3$	$y_4$	$\cdots$	$y_{m-2}$	$y_{m-1}$	$y_m$					
第3趟	$Y$	<table><tr><td><math>y_1</math></td><td><math>y_2</math></td><td><math>y_3</math></td><td><math>y_4</math></td><td><math>\cdots</math></td><td><math>y_{m-2}</math></td><td><math>y_{m-1}</math></td><td><math>y_m</math></td></tr></table>	$y_1$	$y_2$	$y_3$	$y_4$	$\cdots$	$y_{m-2}$	$y_{m-1}$	$y_m$		
$y_1$	$y_2$	$y_3$	$y_4$	$\cdots$	$y_{m-2}$	$y_{m-1}$	$y_m$					
$\cdots$		$\cdots$										

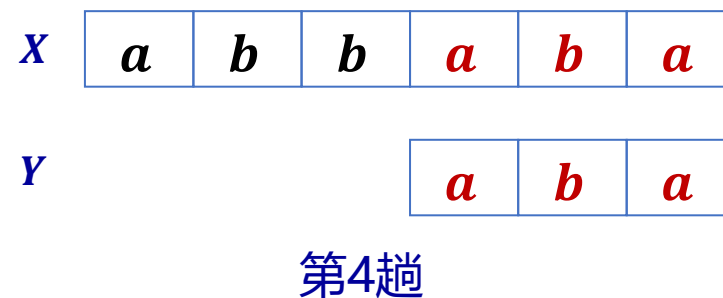
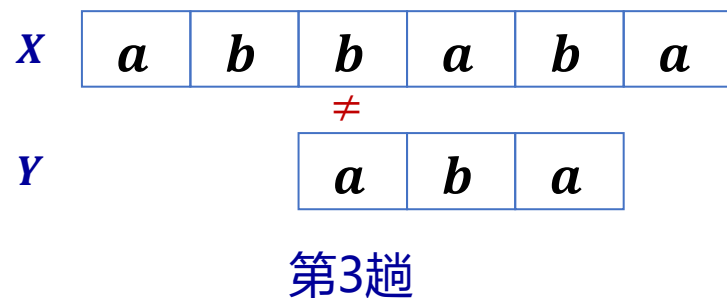
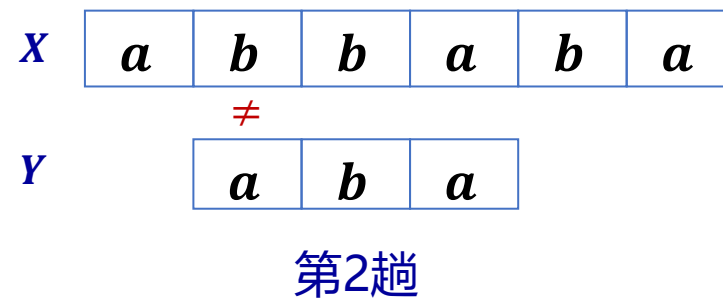
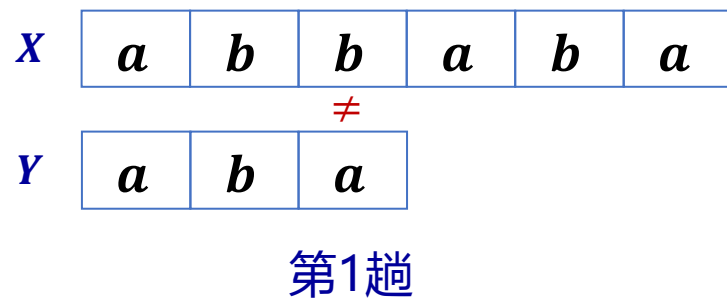
 时间复杂度  $O(mn)$



# 子串匹配问题

- 穷举模式匹配算法

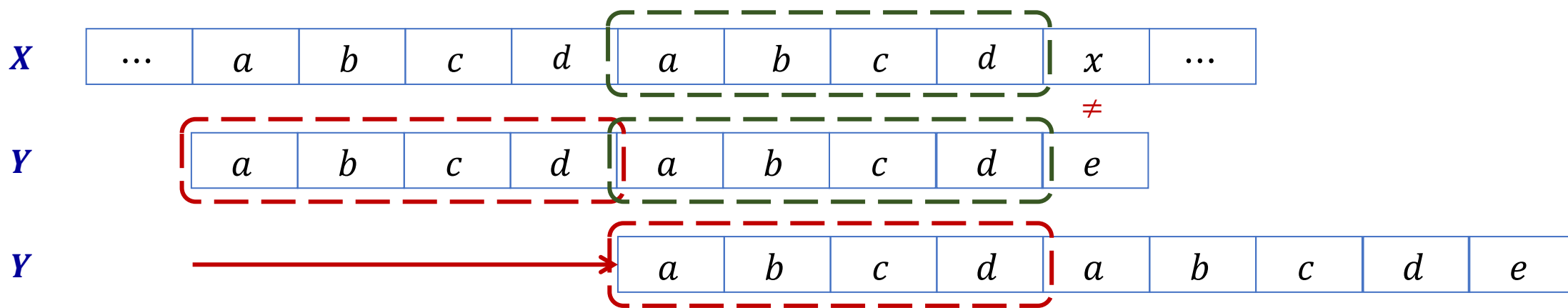
- 匹配失败，目标串 $X$ 回溯，模式串 $Y$ 从头开始





# 子串匹配问题

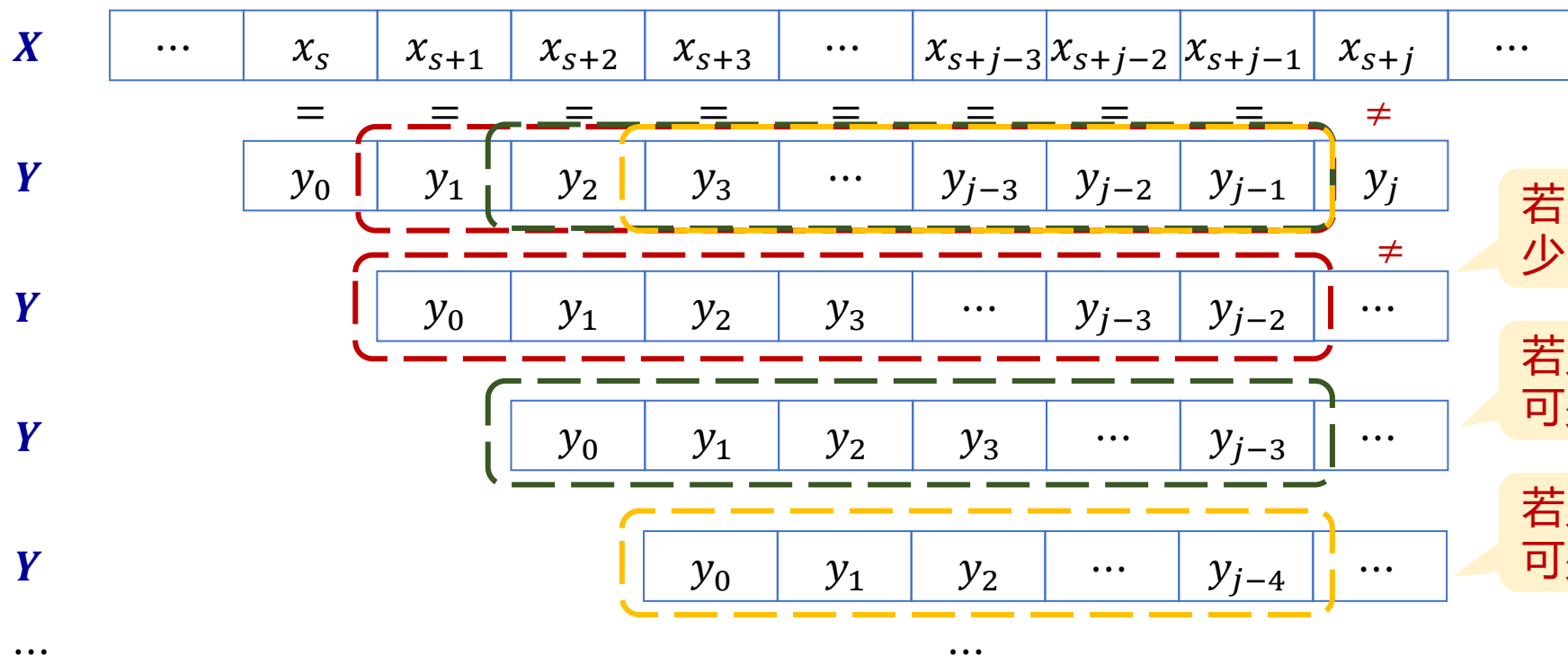
- 确定算法2——Knuth-Morris-Pratt (KMP) 算法
  - 【例】模式串 $Y$ 中重复出现 $abcd$ ，但是 $e$ 和 $x$ 不匹配时，可直接向右滑动4个字符开始匹配，可少匹配3趟





# 子串匹配问题

- Knuth-Morris-Pratt (KMP) 算法



设  $X[s, s + j - 1] = Y[0, j - 1]$ ,  
但  $X[s, s + j] \neq Y[0, j]$

若  $Y[0, j - 2] \neq Y[1, j - 1]$ , 可  
少匹配1趟

若又  $Y[0, j - 3] \neq Y[2, j - 1]$ ,  
可少匹配2趟

若又  $Y[0, j - 4] \neq Y[3, j - 1]$ ,  
可少匹配3趟

类推直到前缀  $Y[0, k + 1] \neq$  后缀  $Y[j - k - 2, j - 1]$ ,  
但是前缀  $Y[0, k] =$  后缀  $Y[j - k - 1, j - 1]$  时,  
可少匹配  $j - k - 2$  趟, 相当于  $Y$  直接向右滑动  $j - k - 1$  个字符





# 子串匹配问题

- Knuth-Morris-Pratt (KMP) 算法
  - 对模式串 $Y$ 进行预处理, 计算可以滑过多少个字符

$$next(j) = \begin{cases} -1, & \text{当 } j = 0 \\ k + 1, & \text{当 } 0 \leq k < j - 1, \text{ 且使 } y_0 y_1 \cdots y_k = y_{j-k-1} y_{j-k} \cdots y_{j-1} \text{ 的最大数} \\ 0, & \text{其他情况} \end{cases}$$

$next(j)$ 直观含义:  $[0, j - 1]$ 中前缀和后缀相等的最大长度

$next(j)$ 直观作用: 可滑过 $j - next(j)$ 位不用匹配

下标 $j$	0	1	2	3	4	5	6	7
$Y$	$a$	$b$	$c$	$d$	$a$	$b$	$c$	$e$
$next(j)$	-1	0	0	0	0	1	2	3



# 子串匹配问题

下标 $j$	0	1	2	3	4	5	6	7
$Y$	$a$	$b$	$c$	$d$	$a$	$b$	$c$	$e$
$next(j)$	-1	0	0	0	0	1	2	3

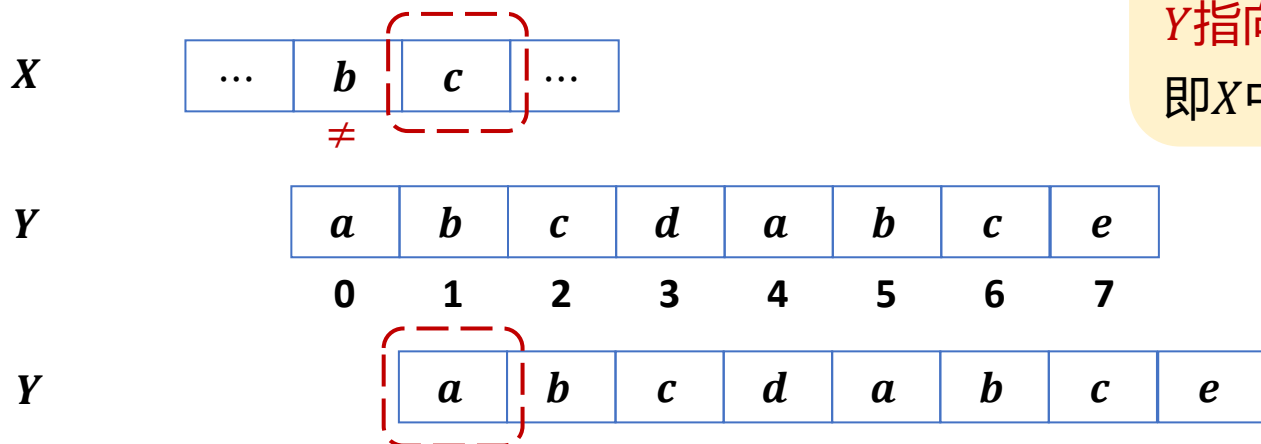
- Knuth-Morris-Pratt (KMP) 算法

- 匹配过程：遇到 $Y[j] \neq X[k]$ 时， $Y$ 向右滑动 $j - next(j)$ 位

$$next(j) = \begin{cases} -1, & \text{当 } j = 0 \\ k + 1, & \text{当 } 0 \leq k < j - 1, \text{ 且使 } y_0 y_1 \cdots y_k = y_{j-k-1} y_{j-k} \cdots y_{j-1} \text{ 的最大数} \\ 0, & \text{其他情况} \end{cases}$$

$j = 0$ 时匹配失败，查找模式串 $Y$ 的 $next(0) = -1$ ，模式串 $Y$ 向右滑动 $j - next(j) = 1$ 位， $X$ 指针加1， $Y$ 指向 $Y[0]$ ，

即 $X$ 中 $c$ 与 $Y$ 中 $Y[0] = a$ 进行比较







# 子串匹配问题

下标 $j$	0	1	2	3	4	5	6	7
$Y$	$a$	$b$	$c$	$d$	$a$	$b$	$c$	$e$
$next(j)$	-1	0	0	0	0	1	2	3

- Knuth-Morris-Pratt (KMP) 算法

- 匹配过程：遇到  $Y[j] \neq X[k]$  时， $Y$  向右滑动  $j - next(j)$  位

$$next(j) = \begin{cases} -1, & \text{当 } j = 0 \\ k + 1, & \text{当 } 0 \leq k < j - 1, \text{ 且使 } y_0 y_1 \cdots y_k = y_{j-k-1} y_{j-k} \cdots y_{j-1} \text{ 的最大数} \\ 0, & \text{其他情况} \end{cases}$$

$X$	...	$a$	$b$	$c$	$d$	$a$	$b$	$c$	$d$	$x$	...
									$\neq$		
$Y$		$a$	$b$	$c$	$d$	$a$	$b$	$c$	$e$		
		0	1	2	3	4	5	6	7		
$Y$						$a$	$b$	$c$	$d$	$a$	$b$
						0	1	2	3	4	5

$j = 7$  时匹配失败，查找模式串  $Y$  的  $next(7) = 3$ ，  
 $X$  中的  $d$  直接与  $Y[next(7)] = Y[3] = d$  进行比较，  
即模式串  $Y$  向右滑动  $j - next(7) = 7 - 3 = 4$  位



# 子串匹配问题

下标 $j$	0	1	2	3	4	5	6	7
$Y$	$a$	$b$	$c$	$d$	$a$	$b$	$c$	$e$
$next(j)$	-1	0	0	0	0	1	2	3

- Knuth-Morris-Pratt (KMP) 算法

- 匹配过程：遇到 $Y[j] \neq X[k]$ 时， $Y$ 向右滑动 $j - next(j)$ 位

$$next(j) = \begin{cases} -1, & \text{当 } j = 0 \\ k + 1, & \text{当 } 0 \leq k < j - 1, \text{ 且使 } y_0 y_1 \cdots y_k = y_{j-k-1} y_{j-k} \cdots y_{j-1} \text{ 的最大数} \\ \mathbf{0}, & \text{其他情况} \end{cases}$$

$X$ 

...	$a$	$b$	$c$	$d$	$a$	$b$	$c$	$d$	$x$	...
-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----

$\neq$

$Y$

$a$	$b$	$c$	$d$	$a$	$b$	$c$	$e$
0	1	2	3	4	5	6	7

$Y$

$a$	$b$	$c$	$d$	$a$	$b$	$c$	$d$	$e$
0	1	2	3	4	5	6	7	8

$j = 4$ 时匹配失败，查找模式串 $Y$ 的 $next(4) = 0$ ， $X$ 中的 $x$ 直接与 $Y[next(4)] = Y[0] = a$ 进行比较，即模式串 $Y$ 向右滑动 $j - next(4) = 4 - 0 = 4$ 位



# 子串匹配问题

- Knuth-Morris-Pratt (KMP) 算法

KMP ( $X, Y, next$ )

$j = 0;$

**for**  $i = 0$  **to**  $n - 1$  **do**

**while**  $j > 0$  and  $Y[j] \neq X[i]$  **do**

$j = next(j);$      /\*向右滑动 $j - next(j)$ 位\*/

**if**  $Y[j] == X[i]$  **then**

$j = j + 1;$

**if**  $j == m$  **then**     /\*找到一个匹配子串\*/

**Print** "Partten occurs with shift"  $i - m;$

$j = next(j);$      /\*继续找下一个匹配字符\*/

- 如果已经预处理得到所有 $next(j)$

 时间复杂度 $\theta(n)$

 如何求所有的 $next(j)$  ?



# 子串匹配问题

下标 $j$	0	1	2	3	4	5	6	7
$Y$	$a$	$b$	$c$	$d$	$a$	$b$	$c$	$d$
$next(j)$	-1	0	0	0	0	1	2	3

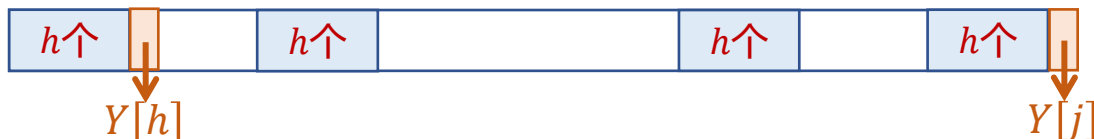
## 快速求 $next(j)$ 的方法

### 假设已知 $next(j) = k$ , 现在计算 $next(j + 1)$

- 若 $Y[k] = Y[j]$ , 则 $next(j + 1) = k + 1 = next(j) + 1$



- 否则, 设 $next(k) = h$ , 此时有 $Y[0, h - 1] = Y[k - h, k - 1] = Y[j - h, j - 1]$



- 若 $Y[h] = Y[j]$ , 则 $next(j + 1) = h + 1$
- 否则, 令 $next(h) = t$ , 此时有 $Y[0, t - 1] = Y[h - t, h - 1] = Y[j - t, j - 1]$ 
  - 继续判断是否 $Y[t] = Y[j]$ , 直到找到或者到 $next(0) = -1$



# 子串匹配问题

- Knuth-Morris-Pratt (KMP) 算法

## Next-Function (Y)

$j = 0; k = -1; \text{next}(0) = -1;$

**while**  $j < m$  **do**

**if**  $k == -1$  or  $Y[j] == Y[k]$  **then**

$j = j + 1;$

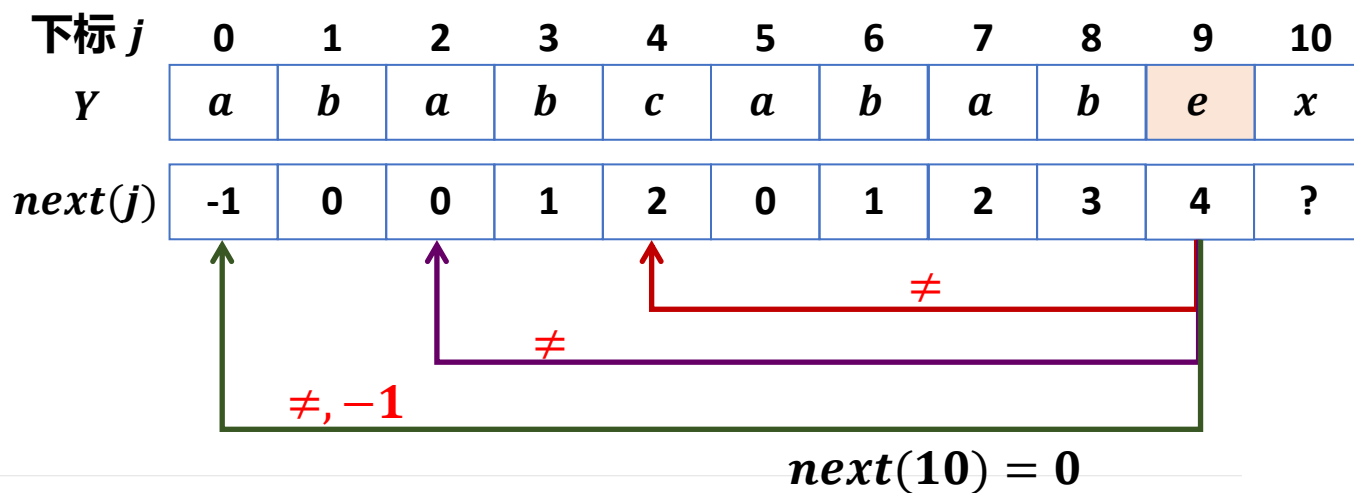
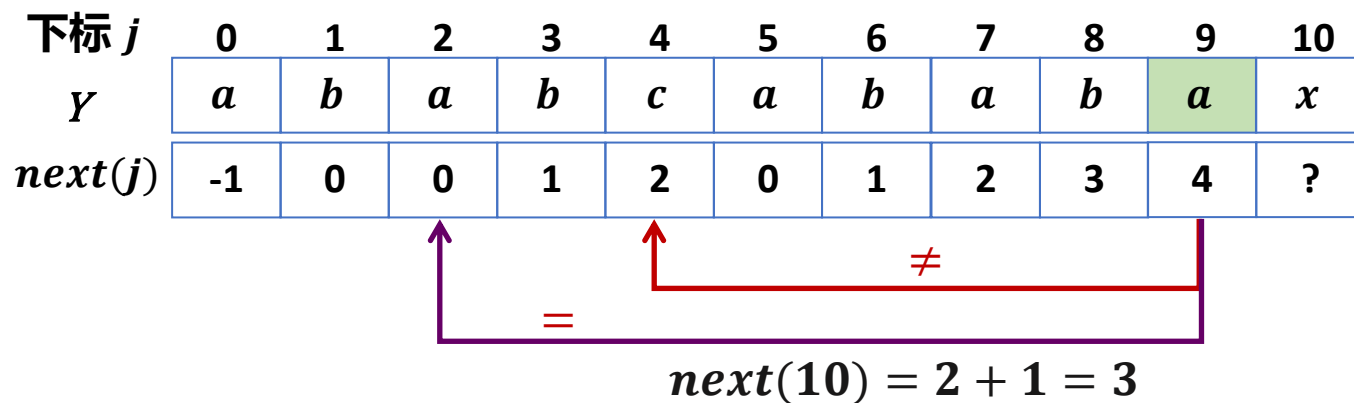
$k = k + 1;$

$\text{next}(j) = k;$

**else**

$k = \text{next}(k);$

**return**  $\text{next};$



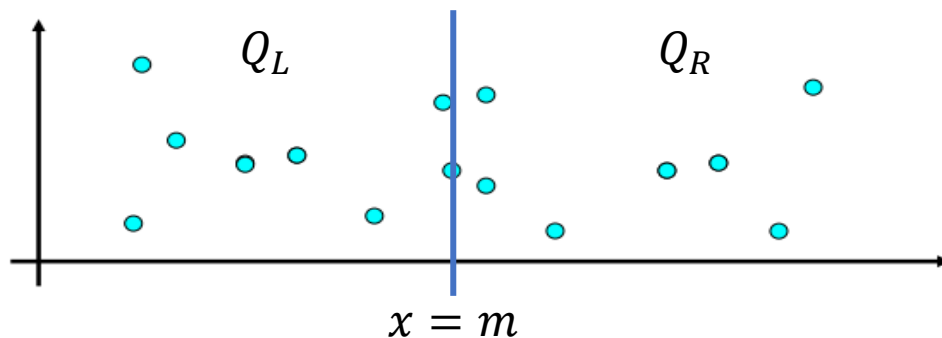


# 最近点对随机算法

- 问题定义

- 输入：欧几里得空间上 $n$ 个点的集合 $Q$
- 输出： $A, B \in Q$ ,  $dis(A, B) = \min\{dis(P_i, P_j) | P_i, P_j \in Q\}$

- 分治算法求解需 $O(n \log n)$ 时间



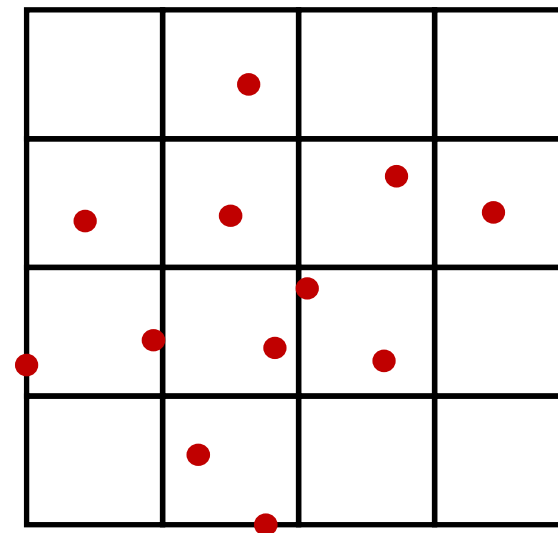
- 用随机算法求解可达到 $O(n)$ 期望时间



# 最近点对随机算法

- 随机算法

- 平面上的 $n$ 个点 $(x_1, x_2)$ ,  $(i = 1, 2, \dots, n)$ , 以 $\Delta$ 为尺寸构造网格,  $\Delta \geq$ 最近距离 $\Delta^*$ 
  - 网格以 $(\min\{x_i\}, \min\{y_j\})$ 为最左下点
  - 以步长 $\Delta$ 不断向右向上伸展
  - $(\max\{x_i\}, \max\{y_j\})$ 含在某一格

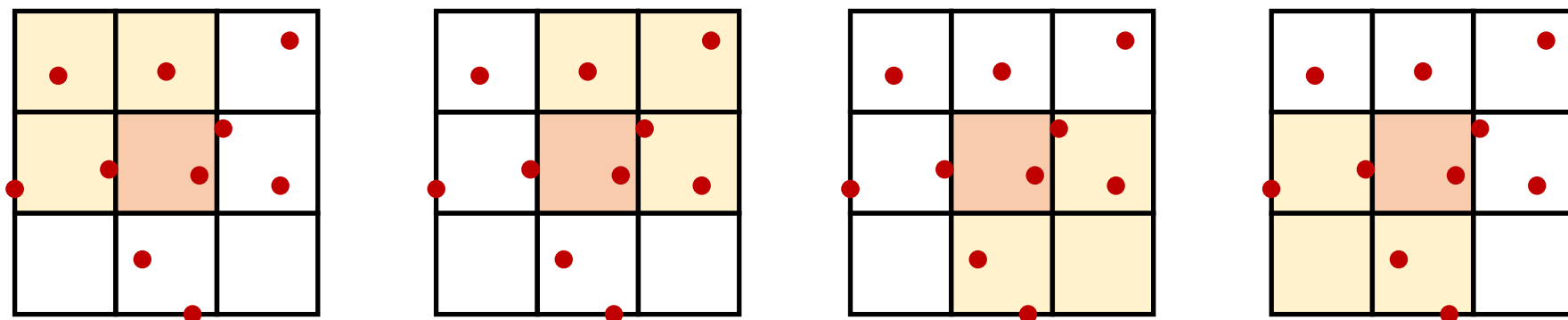




# 最近点对随机算法

- 随机算法

- 平面上的 $n$ 个点 $(x_1, x_2)$ ,  $(i = 1, 2, \dots, n)$ , 以 $\Delta$ 为尺寸构造网格,  $\Delta \geq$ 最近距离 $\Delta^*$ 
  - 若最近点对中的一个点落在某个 $\Delta \times \Delta$ 方格 $C_\Delta$ 中, 则包含 $C_\Delta$ 的4个 $2\Delta \times 2\Delta$ 方格 $C_\Delta^i$  ( $i = 1, 2, 3, 4$ )中, 至少有一个同时含有最近点对的两个点。



尽量给定更小的 $\Delta$





# 最近点对随机算法

## • 随机算法

- 平面上的 $n$ 个点 $(x_1, x_2)$ ,  $(i = 1, 2, \dots, n)$ , 以 $\Delta$ 为尺寸构造网格,  $\Delta \geq$ 最近距离 $\Delta^*$ 。 **如何选择 $\Delta$ ?**

$O(n)$  1. 点集 $S(|S| = n)$ 中随机取子集 $T$ , 使得 $|T| = \lfloor \sqrt{n} \rfloor$

$O(n)$  2.  $T$ 中点两两计算距离, 求出 $T$ 中最近点对距离 $\Delta(T)$

3. 以 $\Delta(T)$ 为尺寸构造网格, 设网格横向 $m$ 个纵向 $r$ 个

4. 分情况讨论:

$O(n)$  A. 如果 $m \times r \leq cn$ , 其中 $c$ 为2~10之间的常数: 找到具有最多节点的方格, 坐标记为 $(i, j)$ ,

$O(n)$  a. 如果方格 $(i, j)$ 包含的节点数 $\leq \sqrt{n}$ , 逐一计算方格 $(i, j)$ 内任意点对的距离, 找到最小值记为 $\Delta'$

期望 $O(n)$  b. 如果方格 $(i, j)$ 包含的节点数 $> \sqrt{n}$ , 将方格 $(i, j)$ 一分为四, 点数最多的子方格点数若还 $> \sqrt{n}$ , 继续拆分直至点数 $\leq \sqrt{n}$ , 计算节点数最多方格的最近点距 $\Delta'$



# 最近点对随机算法

## • 随机算法

- 平面上的 $n$ 个点 $(x_1, x_2)$ ,  $(i = 1, 2, \dots, n)$ , 以 $\Delta$ 为尺寸构造网格,  $\Delta \geq$ 最近距离 $\Delta^*$ 。 **如何选择 $\Delta$ ?**

### 4. 分情况讨论:

A. 如果 $m \times r \leq cn$ , .....

**$O(n)$**  B. 如果 $m \times r > cn$ , 找一个略小于 $cn$ 的素数 $p$ , 用散列函数 $H(i, j, p)$ 把 $S$ 中的点存至长为 $p$ 的散列表中, 其对应一个方格 $(i, j)$ , 散列表的每个槽只对应一个方格 (若发生冲突, 用冲突处理机制), 在散列表中找到点数最多的方格 $(i, j)$ , 然后按照和4.A中同样的方法找到新的 $\Delta'$

**期望 $O(n)$**  5. 以 $\Delta'$ 为新的尺寸构造新网格, 最近点对必落在某个 $2\Delta' \times 2\Delta'$ 的方格中, 找到所有 $2\Delta' \times 2\Delta'$ 的方格中包含节点数大于2的, 并计算其中点对距离, 找到最小值



# 素数测试

- 非对称密码RSA加密算法用到大素数

公钥: $(n, e)$	$n$ 是两个大素数 $p$ 和 $q$ 的乘积 ( $p$ 和 $q$ 必须保密), $e$ 与 $(p-1)(q-1)$ 互质
私钥: $(n, d)$	$e \times d \equiv 1 \pmod{(p-1)(q-1)}$
加密	$C = M^e \pmod{n}$
解密	$M = C^d \pmod{n}$

【例】 令 $p = 3$ ,  $q = 11$ , 则 $n = p \times q = 3 \times 11 = 33$ ;  $(p-1)(q-1) = 2 \times 10 = 20$ ;  
取 $e = 3$ , ( $e = 3$ 与 $(p-1)(q-1) = 20$ 互质)

$3 \times d \equiv 1 \pmod{20}$ , 可取 $d = 7$

明文25, 加密得密文 $16 = 25^3 \pmod{33}$

密文16, 解密得明文 $25 = 16^7 \pmod{33}$

其中关键过程是得到两个大素数,  
但是判断一个非常大的数是否是  
素数并非简单的事



# 素数测试

- 大整数的素因子分解
  - 如 $n = 10^{60}$  (比已知最大素数小很多)
  - 若算法时间复杂度为 $O(n)$
  - 设高速计算机, 每秒1亿亿次( $10^{16}$ )基本运算
  - 则需 $10^{44}$ 秒, 大于 $10^{42}$ 分钟, 大于 $10^{40}$ 小时, 大于 $10^{38}$ 天, 大于 $10^{35}$ 年!
  - $O(n)$ 时间的算法绝对不能接受!



# 素数测试

- 准备工作——求 $a^m \pmod n$ 的算法 ( $m \leq n$ )
  - $m$ 的二进制表示为 $b_k b_{k-1} \cdots b_1 b_0$ , ( $b_k = 1, k \approx \log_2 m$ )
  - $a^m = a^{b_k b_{k-1} \cdots b_1 b_0}$
  - 【例】  $m = 41 = b_5 b_4 b_3 b_2 b_1 b_0 = 101001_{(2)}$ ,  $k = 5$
  - $a^{41} = a^{101001_{(2)}}$



# 素数测试

$m$  的二进制表示为  $b_k b_{k-1} \cdots b_1 b_0$ , ( $b_k = 1, k \approx \log_2 m$ )

- 准备工作——求  $a^m \pmod n$  的算法 ( $m \leq n$ )
  - 计算  $a^m$  可以用下述方法
    - 从  $m$  的二进制的高位到低位, 平方, 遇1还要乘  $a$
  - 【例】计算  $a^{41}$ ,  $41 = b_5 b_4 b_3 b_2 b_1 b_0 = 101001_{(2)}$ 
    - 初始  $C = 1$
    - $b_5 = 1$ :  $C = C^2 = 1$ , 因为  $b_5 = 1$ ,  $C = a * C = a$
    - $b_5 b_4 = 10$ :  $C = C^2 = a^2$
    - $b_5 b_4 b_3 = 101$ :  $C = C^2 = a^4$ , 因为  $b_3 = 1$ ,  $C = a * C = a^5$
    - $b_5 b_4 b_3 b_2 = 1010$ :  $C = C^2 = a^{10}$
    - $b_5 b_4 b_3 b_2 b_1 = 10100$ :  $C = C^2 = a^{20}$
    - $b_5 b_4 b_3 b_2 b_1 b_0 = 101001$ :  $C = C^2 = a^{40}$ , 因为  $b_0 = 1$ ,  $C = a * C = a^{41}$



# 素数测试

$m$  的二进制表示为  $b_k b_{k-1} \cdots b_1 b_0$ , ( $b_k = 1, k \approx \log_2 m$ )

## • 准备工作——求 $a^m \pmod n$ 的算法 ( $m \leq n$ )

### • 计算 $a^m$ 可以用下述方法

- 从  $m$  的二进制的高位到低位, 平方, 遇1还要乘  $a$

好处:  $C$  的最大值不超过  $n-1$ , 中间值不超过  $\max\{(n-1)^2, a(n-1)\}$  求  $a^m \pmod n$  时不会占用很多空间

### • 模运算有规则: $(x * y) \% n = ((x \% n) * (y \% n)) \% n$

### • 求 $a^m \pmod n$ 可在求 $a^m$ 的每一步求模

**EXPMOD** ( $a, m, n$ )

$C = 1;$

**for**  $j = k$  **to** 0

$C = C^2 \pmod n;$

**if**  $b_j = 1$  **then**

$C = a * C \pmod n;$

**return**  $C;$

例如, 计算  $a^5 \pmod n$ ,  $m = 5 = 101_{(2)}$

- 初始  $C = 1$
- $b_2 = 1$ :  $C = C^2 \pmod n$ , 因为  $b_2 = 1$ ,  $C = a * C \pmod n$
- $b_2 b_1 = 10$ :  $C = C^2 \pmod n$
- $b_2 b_1 b_0 = 101$ :  $C = C^2 \pmod n$ , 因为  $b_0 = 1$ ,  $C = a * C \pmod n$

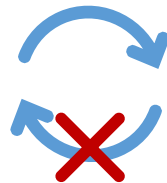


# 素数测试

- Fermat小定理

- 若 $n$ 为素数,  $0 < a < n$ , 则有 $a^{n-1} \equiv 1 \pmod{n}$ 
  - 可以从 $a = 2, 3, \dots, n-1$ 对 $a^{n-1} \equiv 1 \pmod{n}$ 进行测试, 如果一旦出现不满足, 则 $n$ 不是素数
  - 但是, 如果全部满足 $a^{n-1} \equiv 1 \pmod{n}$ 也不能说明 $n$ 是素数
  - 某些 $n$ 是合数, 对任意 $a$ , 也满足 $a^{n-1} \equiv 1 \pmod{n}$ , 称之为Carmichael数, 例如前几个561, 1105, 1729, 2465, ..., 小于1亿只有255个Carmichael数
  - Carmichael数虽然分布很稀, 但仍有无穷多个

$n$ 为素数,  $0 < a < n$



$$a^{n-1} \equiv 1 \pmod{n}$$





# 素数测试

- Miller-Rabin算法——剔除Carmichael数
  - $1 < a < n$ , 有  $a^{n-1} \equiv 1 \pmod{n}$ 
    - 前提:  $n$  为奇数 (如果为偶数就不用判断它是否为素数了)
    - 那么,  $n-1$  为偶数, 即  $a^{n-1}$  可以转换为  $a^{2k}$ , 即  $a^{2k} \equiv 1 \pmod{n}$
    - 如果  $a^{2k} \equiv 1 \pmod{n}$ , 那么,  $n$  整除  $a^{2k} - 1$
    - $\because a^{2k} - 1 = (a^k)^2 - 1 = (a^k - 1)(a^k + 1)$
    - $\therefore$  如果  $n$  为素数, 那么  $n$  整除  $(a^k - 1)$  或整除  $(a^k + 1)$
    - 那么,  $a^k \pmod{n} = \begin{cases} 1 \\ n-1 \end{cases}$
  - 二次探测: 测试  $a^{n-1} \equiv 1 \pmod{n}$  之后, 继续测试  $a^k \pmod{n}$ , 其中,  $2k = n-1$



# 素数测试

- Miller-Rabin算法

- $1 < a < n$ , 有  $a^{n-1} \equiv 1 \pmod{n}$
- 二次探测: 测试  $a^{n-1} \equiv 1 \pmod{n}$  之后, 继续测试  $a^k \pmod{n} = 1$  或  $n-1$ , 其中,  $2k = n-1$
- 继续探测:
  - 找到  $n-1 = m \times 2^q$ , 其中,  $m$  为奇数。
  - 用  $a^m, a^{m \times 2}, a^{m \times 2^2}, a^{m \times 2^3}, \dots, a^{m \times 2^q} = a^{n-1}$  分别对  $n$  取模, 一旦结果不为 1 或  $n-1$ , 则判断  $n$  为合数, 此时的  $a$  称为证据数 (witness)



在  $1 \sim n$  之间的证据数极多, 但目前理论上只能证明它们超过  $(n-1)/2$ , 即至少有一半以上的证据数。于是答错的概率最多为  $1/2$ , (实际上低得多)。随机选  $k$  个证据数  $a$ , 答错概率小于  $1/2^k$



# 素数测试

## • Miller-Rabin算法

- 定理：设 $n$ 为素数，且 $0 < x < n$ ，则当 $x^2 \equiv 1 \pmod{n}$ 时，必有 $x = 1$ 或 $x = n - 1$
- 推论：当 $0 < x < n$ 且 $x^2 \equiv 1 \pmod{n}$ 成立时，若 $x \neq 1$ 且 $x \neq n - 1$ ，则 $n$ 是合数

```
Miller-Rabin ( $m, n$ )  
 $a = \text{random}(1, n - 1);$   
 $x = a^m \pmod{n};$   
for  $i = 1$  to  $q$  do  
     $y = x^2 \pmod{n};$   
    if  $y = 1$  and  $x \neq 1$  and  $x \neq n - 1$  then /*此时满足推论*/  
        return “ $n$ 为合数”;  
     $x = y;$   
if  $x \neq 1$  then return “ $n$ 为合数”;  
else  
    return “ $n$ 为素数”; /*此时满足Fermat小定理*/
```



# 素数测试

- 【例】测试Carmichael数561是否为素数

- $n - 1 = 560 = 2^4 * 35$
- 即有  $q = 4, m = 35$
- 假定选  $a = 7$ , 则  $a^m = 7^{35} \equiv 241 \pmod{561}$
- $7^{35}, 7^{70}, 7^{140}, 7^{280}, 7^{560} \pmod n$  后分别为: 241, 298, 166, 67, 1
- $x = 7^{280} \equiv 67 \pmod{561}$ ,
- 而  $x^2 = 7^{560} \equiv 1 \pmod{561}$ ,  $x \neq 1$  且  $x \neq n - 1$
- 则561是合数

## Miller-Rabin ( $m, n$ )

$a = \text{random}(1, n - 1);$

$x = a^m \pmod n;$

**for**  $i = 1$  **to**  $q$  **do**

$y = x^2 \pmod n;$

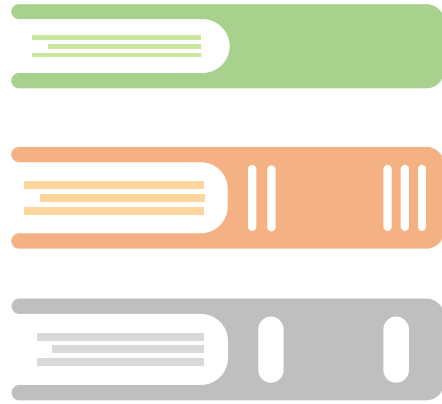
**if**  $y = 1$  and  $x \neq 1$  and  $x \neq n - 1$  **then**  
**return** “ $n$ 为合数”;

$x = y;$

**if**  $x \neq 1$  **then return** “ $n$ 为合数”;

**else**

**return** “ $n$ 为素数”;



**Q&A**