

# random variables

- A student is taking a true-false test and guessing when he doesn't know the answer. We are going to compute a score by subtracting a percentage of the number of incorrect answers from the number of correct answers. That is, for some number  $y$ , the student's corrected score will be

$$(\text{number of corrected answers}) - y(\text{number of incorrect answers})$$

When we convert this “corrected score” to a percentage score, we want its expected value to be the percentage of the material being tested that the student knows. How can we do this?

# 计算机问题求解 – 论题2-6

- 概率分析与随机算法

课程研讨

- TC第5章
- CS第5章第6、7节

# 问题1： randomized algorithm

- 什么样的算法可以称作randomized algorithm?
- 什么叫做randomized algorithm的expected running time?
- 它和average-case running time有什么异同?

# 问题1： randomized algorithm

- 什么样的算法可以称作randomized algorithm?
  - Its behavior is determined not only by its input but also by something chosen randomly (e.g. values produced by a random-number generator).
- 什么叫做randomized algorithm的expected running time?
- 它和average-case running time有什么异同?

# 问题1： randomized algorithm

- 什么样的算法可以称作randomized algorithm?
  - Its behavior is determined not only by its input but also by something chosen randomly (e.g. values produced by a random-number generator).
- 什么叫做randomized algorithm的expected running time?
- 它和average-case running time有什么异同?
  - 异： We discuss the average-case running time when the probability distribution is over the inputs to the algorithm, and we discuss the expected running time when the algorithm itself makes random choices.

# 问题1： randomized algorithm (续)

- 你能想到哪些方法生成一个32-bit的（伪）随机数？

# 问题1： randomized algorithm (续)

- 你能想到哪些方法生成一个32-bit的（伪）随机数？
  - Computational methods (pseudo-random number generators)

```
m_w = <choose-initializer>; /* must not be zero */
m_z = <choose-initializer>; /* must not be zero */

uint get_random()
{
    m_z = 36969 * (m_z & 65535) + (m_z >> 16);
    m_w = 18000 * (m_w & 65535) + (m_w >> 16);
    return (m_z << 16) + m_w; /* 32-bit result */
}
```

- Physical methods
  - Coin flipping
  - Dice
  - Variations in the amplitude of atmospheric noise recorded with a normal radio

# 问题1： randomized algorithm (续)

- 你能想到哪些方法对一个数组中的元素随机排序？  
你如何评价这样一个方法的好坏？



# 问题1: randomized algorithm (续)

- 你能想到哪些方法对一个数组中的元素随机排序? 你如何评价这样一个方法的好坏?

- PERMUTE-BY-SORTING( $A$ )

```
1  $n = A.length$   
2 let  $P[1..n]$  be a new array  
3 for  $i = 1$  to  $n$   
4      $P[i] = \text{RANDOM}(1, n^3)$   
5 sort  $A$ , using  $P$  as sort keys
```

- RANDOMIZE-IN-PLACE( $A$ )

```
1  $n = A.length$   
2 for  $i = 1$  to  $n$   
3     swap  $A[i]$  with  $A[\text{RANDOM}(i, n)]$ 
```

## 问题2: expected running time

- 目前为止, 你掌握了哪些方式计算 $E(X)$ ?

## 问题2: expected running time

- 目前为止, 你掌握了哪些方式计算 $E(X)$ ?
  - $E(X) = \sum xP(X=x)$  // 定义
  - $E(X) = \sum E(X_i)$  // indicator random variable
  - $E(aX+bY) = aE(X) + bE(Y)$  // linearity of expectation
  - $E(X) = \sum E(X|F_i)P(F_i)$  // conditional expected value

## 问题2: expected running time (续)

- 你能解释这里用的是哪种方法吗?

Slower Quicksort(A,n)

if ( $n = 1$ )

    return the one item in  $A$

else

    Repeat

$p = \text{randomElement}(A)$

        Let  $H$  be the set of elements greater than  $p$ ; Let  $h = |H|$

        Let  $L$  be the set of elements less than or equal to  $p$ ; Let  $\ell = |L|$

    Until ( $|H| \geq n/4$ ) and ( $|L| \geq n/4$ )

$A_1 = \text{QuickSort}(H, h)$

$A_2 = \text{QuickSort}(L, \ell)$

    return the concatenation of  $A_1$  and  $A_2$

$$T(n) \leq E(r)bn + T(a_n n) + T((1 - a_n)n)$$

## 问题2: expected running time (续)

- 你能解释这里用的是哪种方法吗?

RandomSelect(A,i,n)

(selects the  $i$ th smallest element in set  $A$ , where  $n = |A|$  )

if ( $n = 1$ )

    return the one item in  $A$

else

$p = \text{randomElement}(A)$

    Let  $H$  be the set of elements greater than  $p$

    Let  $L$  be the set of elements less than or equal to  $p$

    If ( $H$  is empty)

        put  $p$  in  $H$

    if ( $i \leq |L|$ )

        Return RandomSelect( $L, i, |L|$ )

    else

        Return RandomSelect( $H, i - |L|, |H|$ ).

$$T(n) \leq \begin{cases} \frac{1}{2}T(\frac{3}{4}n) + \frac{1}{2}T(n) + bn & \text{if } n > 1 \\ d & \text{if } n = 1 \end{cases}$$

## 问题2: expected running time (续)

- 你能解释这里用的是哪种方法吗?

**Exercise 5.6-4** Consider an algorithm that, given a list of  $n$  numbers, prints them all out. Then it picks a random integer between 1 and 3. If the number is 1 or 2, it stops. If the number is 3 it starts again from the beginning. What is the expected running time of this algorithm?

$$T(n) = \frac{2}{3}cn + \frac{1}{3}(cn + T(n))$$

## 问题2： expected running time (续)

- 你怎么理解indicator random variable?
- 如何利用indicator random variable来简化期望的计算?
  
- 在这些问题中， indicator random variable分别可以是什么？
  - The expected number of times that we hire a new office assistant.
  - The expected number of pairs of people with the same birthday.
  - How many sixes do we expect to see on top if we roll 24 dice?  
(上周，根据期望的定义，我们是如何计算的？)

## 问题2: expected running time (续)

- 你怎么理解indicator random variable?
- 怎么利用indicator random variable来简化期望的计算?

$$\begin{aligned} E[X] &= E\left[\sum_{i=1}^n X_i\right] \\ &= \sum_{i=1}^n E[X_i] \end{aligned}$$

Given a sample space  $S$  and an event  $A$  in the sample space  $S$ , let  $X_A = I\{A\}$ .  
Then  $E[X_A] = \Pr\{A\}$ .

- 在这些问题中, indicator random variable分别可以是什么?
  - The expected number of times that we hire a new office assistant.
  - The expected number of pairs of people with the same birthday.
  - How many sixes do we expect to see on top if we roll 24 dice?  
(上周, 根据期望的定义, 我们是如何计算的?)



## 问题2: expected running time (续)

- Suppose that you want to output 0 with probability  $1/2$  and 1 with probability  $1/2$ . At your disposal is a procedure **BIASED-RANDOM**, that outputs either 0 or 1. It outputs 1 with some probability  $p$  and 0 with probability  $1 - p$ , where  $0 < p < 1$ , but you do not know what  $p$  is. Give an algorithm that uses **BIASED-RANDOM** as a subroutine, and returns an unbiased answer, returning 0 with probability  $1/2$  and 1 with probability  $1/2$ . What is the expected running time of your algorithm as a function of  $p$ ?

## 问题2: expected running time (续)

- UNBIASED-RANDOM()

**Output:** 0 with probability 1/2 and 1 with probability 1/2

```
1 while true do
2   | a ← BIASED-RANDOM()
3   | b ← BIASED-RANDOM()
4   | if a < b then return 0
5   | if a > b then return 1
```

The algorithm calls BIASED-RANDOM twice to get two random numbers  $A$  and  $B$ . It repeats this until  $A \neq B$ . Then, depending on whether  $A < B$  (that is,  $A = 0$  and  $B = 1$ ) or  $A > B$  (that is,  $A = 1$  and  $B = 0$ ) it returns 0 or 1 respectively.

In any iteration, we have  $\Pr(A < B) = p(1 - p) = \Pr(B < A)$ , that is, the probability that the algorithm returns 0 in that iteration equals to the probability that it returns 1 in that iteration. Since with probability 1 we return something at some point (and not repeat the loop endlessly) and the probabilities of returning 0 and 1 are equal in each iteration, the total probabilities of returning 0 and 1 must be 1/2 and 1/2 respectively.

- 怎么计算expected running time?

## 问题2: expected running time (续)

- UNBIASED-RANDOM()

**Output:** 0 with probability  $1/2$  and 1 with probability  $1/2$

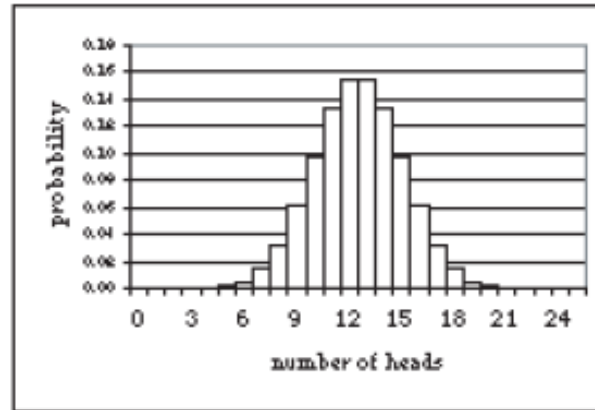
```
1 while true do
2    $a \leftarrow$  BIASED-RANDOM()
3    $b \leftarrow$  BIASED-RANDOM()
4   if  $a < b$  then return 0
5   if  $a > b$  then return 1
```

- 怎么计算expected running time?

The algorithm stops, if it either returns 0 or 1. In every iteration, the probability of this is  $\Pr(A \neq B) = \Pr(A < B) + \Pr(B < A) = 2p(1 - p)$ . Thus, we have a sequence of independent Bernoulli trials, each with probability  $2p(1 - p)$  of success. Therefore, the number of iterations required before the algorithm stops is geometrically distributed with parameter  $2p(1 - p)$ , and the expected number of iterations is  $1/(2p(1 - p))$ . As each iteration takes constant time (assuming that BIASED-RANDOM takes constant time), the expected running time of the algorithm is  $\Theta(1/(p(1 - p)))$ .

## 问题3: probability distributions and variance

- 你怎么理解distribution function和它的histogram?



- 你怎么理解cumulative distribution function?  
它有哪些性质?  
什么情况下只能使用cumulative distribution function?

# 问题3: probability distributions and variance

(续)

- 谈谈你对variance的理解
- 为什么variance被定义成 $E((X-E(X))^2)$ 这种形式?