



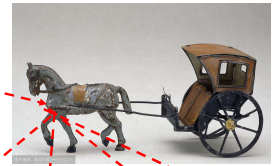
## 第5章 鞅

骆吉洲

计算机科学与技术学院



鞅: martingale



四川老乡读出来:

“鞅起弟弟/妹妹一路走...”

含义: 带弟弟/妹妹一起走, 以他/她的速度来限制你的速度

随机变量序列:  $X_0, X_1, \dots, X_n, \dots$

限制  $X_0, X_1, \dots, X_n, \dots$  增长速度的内在机制



估计硬币投掷过程中

- 正面向上次数与背面向上次数之差的大小



估计随机图 (社交网络、路由网络) 中

- 路径长短
- 围长
- 色数
- 团大小
- ...



估计长序列 (DNA序列、字符串序列) 中

- 模式串出现次数



- 有依赖关系的随机变量序列  $X_0, X_1, X_2, \dots$

$$\Pr[|X_n - X_0| > t] = ?$$

- 随机变量  $X_1, X_2, \dots, X_n$  的函数  $f(X_1, X_2, \dots, X_n)$  不仅仅是和

$$\Pr[|f(X_1, X_2, \dots, X_n) - E[f(X_1, X_2, \dots, X_n)]| > t] = ?$$

哪种结构有助于... 鞅!



### 提纲

#### 5.1 鞅的定义和基本性质

#### 5.2 鞅的一般形式

停时定理 瓦尔德方程 鞅尾不等式

#### 5.3 简单应用

##### 5.3.1 模式匹配

##### 5.3.2 球和箱子模型中空箱子的个数

##### 5.3.3 随机图的色数




### 5.1 鞅的定义和基本性质

HIT CS&E

## 鞅(Martingale)

**定义**  
 随机变量序列  $X_0, X_1, X_2, \dots$   
 如果  $E[X_i | X_0, X_1, \dots, X_{i-1}] = X_{i-1} \quad \forall i \geq 1$   
 则称  $X_0, X_1, X_2, \dots$  是一个鞅



我不知道你在说啥

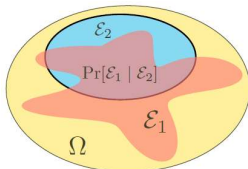
HIT CS&E

## 条件概率

**定义.** “在随机事件  $\mathcal{E}_2$  发生的条件下  $\mathcal{E}_1$  发生” 的条件概率定义为

$$\Pr[\mathcal{E}_1 | \mathcal{E}_2] = \frac{\Pr[\mathcal{E}_1 \wedge \mathcal{E}_2]}{\Pr[\mathcal{E}_2]}$$

如果  $\mathcal{E}_1$  和  $\mathcal{E}_2$  是独立的

$$\begin{aligned} \Pr[\mathcal{E}_1 | \mathcal{E}_2] &= \frac{\Pr[\mathcal{E}_1 \wedge \mathcal{E}_2]}{\Pr[\mathcal{E}_2]} \\ &= \frac{\Pr[\mathcal{E}_1] \Pr[\mathcal{E}_2]}{\Pr[\mathcal{E}_2]} \\ &= \Pr[\mathcal{E}_1] \end{aligned}$$


HIT CS&E

## 条件期望

**定义.**  $Y$  关于  $X$  的条件期望定义为

$$E[Y | X] = \sum y \cdot \Pr[Y=y | X]$$

**例.** 均匀随机地抽取一个人  
 $X$ : 国籍  $Y$ : 身高  $Z$ : 性别

$$E[Y | X = \text{“中国”}] = \sum y \cdot \Pr[Y=y | X = \text{“中国”}]$$

给定一个“国籍”  $x$ , 我们就可以计算  $E[Y | X=x]$   
 $E[Y | X]$ :  $X$  的变化范围  $\rightarrow R$   
 $E[Y | X]$  可以视为一个函数  $f(X)$   
 $f(x) = E[Y | X=x]$   $E[Y | X]$  也是一个随机变量

HIT CS&E

## 条件期望

**例.** 均匀随机地抽取一个人  
 $X$ : 国籍  $Y$ : 身高  $Z$ : 性别

$$E[Y | X, Z]$$

$$E[Y | X = \text{“中国”}, Z = \text{“男”}] = \text{中国男性平均身高}$$

$$E[Y | X = \text{“中国”}, Z = \text{“女”}] = \text{中国女性平均身高}$$

$$E[Y | X = \text{“美国”}, Z = \text{“男”}] = \text{美国男性平均身高}$$

$$E[Y | X = \text{“美国”}, Z = \text{“女”}] = \text{美国女性平均身高}$$

.....  
 $E[Y | X, Z]$  也是一个随机变量

HIT CS&E

## 条件期望的基本性质

**性质1.**  $E[Y] = E[E[Y | X]]$   
**例:** 所有人身高均值 = 各国身高均值的均值

**性质2.**  $E[Y | Z] = E[E[Y | X, Z] | Z]$   
**例:** 所有男人身高均值 = 各国男人身高均值的均值

**性质3.**  $E[E[f(X)g(X, Y) | X]] = E[f(X)] E[g(X, Y) | X]$   
 一旦给定  $X=x$ , 则  $f(X)$  的取值是一个实数

HIT CS&E

## 鞅

**定义**  
 随机变量序列  $X_0, X_1, X_2, \dots$   
 如果  $E[X_i | X_0, X_1, \dots, X_{i-1}] = X_{i-1} \quad \forall i \geq 1$   
 则称  $X_0, X_1, X_2, \dots$  是一个鞅

$\forall X_0, X_1, \dots, X_{i-1},$   
 $E[X_i | X_0 = x_0, X_1 = x_1, \dots, X_{i-1} = x_{i-1}] = x_{i-1}$   
 $E[X_i - X_{i-1} | X_0, \dots, X_{i-1}] = 0$

HIT CS&E

### 例1: 公平赌博

$$E[X_i | X_0, X_1, \dots, X_{i-1}] = X_{i-1}$$

公平赌局

- 初始赌资  $X_0$ ,  $X_i$  表示第  $i$  轮之后的赌资
- 每轮以  $1/2$  的概率赢, 赢得的赌资是  $a$
- 以  $1/2$  的概率输, 输掉的赌资是  $a$
- 每局的输赢是相互独立的

$$E[X_i | X_0, X_1, \dots, X_{i-1}]$$


$$= E[X_i | X_{i-1}]$$

$$= (X_{i-1} + a) \Pr[\text{第 } i \text{ 轮赢}] + (X_{i-1} - a) \Pr[\text{第 } i \text{ 轮输}]$$

$$= (X_{i-1} + a)/2 + (X_{i-1} - a)/2$$

$$= X_{i-1}$$

赌资序列是鞅



HIT CS&E

### 例2: 例1的抽象化

$$E[X_i | X_0, X_1, \dots, X_{i-1}] = X_{i-1}$$

均匀硬币投掷  $n$  次

- $X_0 = 0$
- $X_i = \text{HEADS}_i - \text{TAILS}_i$
- $\text{HEADS}_i$  —  $i$  次投掷后头面向上的总次数
- $\text{TAILS}_i$  —  $i$  次投掷后背面向上的总次数

$$E[X_i | X_0, X_1, \dots, X_{i-1}]$$


$$= E[X_i | X_{i-1}]$$

$$= (X_{i-1} + 1) \Pr[\text{第 } i \text{ 次头面向上}] + (X_{i-1} - 1) \Pr[\text{第 } i \text{ 次背面向上}]$$

$$= (X_{i-1} + 1)/2 + (X_{i-1} - 1)/2$$

$$= X_{i-1}$$

实验成败总次数之差的序列是鞅



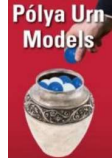
HIT CS&E

### 例3: Polya 壶

$$E[X_i | X_0, X_1, \dots, X_{i-1}] = X_{i-1}$$

Polya 壶

- 壶中有  $b$  个篮球,  $w$  个白球,  $X_0$
- 第  $i$  次操作
  - 均匀、随机、独立地从壶中抓出一个球
  - 放回与所抓球同色的  $c$  个球
  - 第  $i$  次操作后篮球的比例记为  $X_i$
- $X_0, X_1, \dots, X_n, \dots$



HIT CS&E

### 例3: Polya 壶

$$E[X_i | X_0, X_1, \dots, X_{i-1}] = X_{i-1}$$

操作前: 球总数  $a$  篮球比例  $X_{i-1}$

抓到篮球,  $\Pr = X_{i-1}$       抓到白球,  $\Pr = 1 - X_{i-1}$

操作后: 球总数  $a-1+c$       操作后: 球总数  $a-1+c$

篮球个数  $aX_{i-1}-1+c$       篮球个数  $aX_{i-1}$


$$E[X_i | X_0, X_1, \dots, X_{i-1}]$$

$$= E[X_i | X_{i-1}]$$

$$= \frac{aX_{i-1}-1+c}{a-1+c} X_{i-1} + \frac{aX_{i-1}}{a-1+c} (1-X_{i-1})$$

$$= X_{i-1}$$

Polya 壶中篮球比例是鞅



HIT CS&E

### 鞅尾不等式

给定鞅  $X_0, X_1, \dots, X_n, \dots$

$\Pr[|X_n - X_0| > t] = ?$

- 鞅中随机变量偏离  $X_0$  超过阈值的概率
- 鞅中随机变量偏离数学期望超过阈值的概率

限制速度  
火车跑得快, 全靠车头带

HIT CS&E

### Azuma 不等式

Azuma 不等式

如果鞅  $X_0, X_1, X_2, \dots$  对  $k \geq 1$  满足

$$|X_k - X_{k-1}| \leq c_k$$

则

$$\Pr[|X_n - X_0| \geq t] \leq 2 \exp\left(-\frac{t^2}{2 \sum_{k=1}^n c_k^2}\right)$$

对于随机变量序列, 如果每步

- 从平均看, 不会偏离当前的值(鞅)
- 取值不会有大的跳跃

则其最终取值不会偏离初始值太远



## Azuma不等式

### Azuma不等式

如果鞅 $X_0, X_1, X_2, \dots$ 对 $k \geq 1$ 满足

$$|X_k - X_{k-1}| \leq c_k$$

则

$$\Pr[|X_n - X_0| \geq t] \leq 2 \exp\left(-\frac{t^2}{2 \sum_{k=1}^n c_k^2}\right)$$

**推论:** 如果鞅 $X_0, X_1, X_2, \dots$ 对 $k \geq 1$ 满足

$$|X_k - X_{k-1}| \leq c$$

则

$$\Pr[|X_n - X_0| \geq ct\sqrt{n}] \leq 2e^{-t^2/2}$$



## Azuma不等式证明思路

### Azuma不等式

如果鞅 $X_0, X_1, X_2, \dots$ 对 $k \geq 1$ 满足

$$|X_k - X_{k-1}| \leq c_k$$

则

$$\Pr[|X_n - X_0| \geq t] \leq 2 \exp\left(-\frac{t^2}{2 \sum_{k=1}^n c_k^2}\right)$$

**第1步:** 将总差表示为分步差之和

$$Y_i = X_i - X_{i-1} \quad X_n - X_0 = \sum_{i=1}^n Y_i$$

**第2步:** 将markov不等式应用到矩生成函数

$$\Pr[\sum_{i=1}^n Y_i \geq t] = \Pr[e^{\lambda \sum_{i=1}^n Y_i} \geq e^{\lambda t}] \leq \frac{\mathbb{E}[e^{\lambda \sum_{i=1}^n Y_i}]}{e^{\lambda t}}$$

**第3步:** 利用鞅的性质和矩生成函数的凸性质



**推论:** 如果鞅 $X_0, X_1, X_2, \dots$ 对 $k \geq 1$ 满足

$$|X_k - X_{k-1}| \leq c$$

则

$$\Pr[|X_n - X_0| \geq ct\sqrt{n}] \leq 2e^{-t^2/2}$$

**公平赌局**

- 初始赌资 $X_0$ ,  $X_i$ 表示第 $i$ 轮之后的赌资
- 每轮以1/2的概率赢, 赢得的赌资是 $a$
- 以1/2的概率输, 输掉的赌资是 $a$
- 每局的输赢是相互独立的

$$|X_i - X_{i-1}| \leq a \quad \Pr[|X_n - X_0| \geq atn^{1/2}] \leq 2 \exp(-t^2/2)$$

公平赌局中, 有很大的输赢是不太可能的!



**推论:** 如果鞅 $X_0, X_1, X_2, \dots$ 对 $k \geq 1$ 满足

$$|X_k - X_{k-1}| \leq c$$

则

$$\Pr[|X_n - X_0| \geq ct\sqrt{n}] \leq 2e^{-t^2/2}$$

**均匀硬币投掷 $n$ 次**

- $X_0 = 0$
- $X_i = \text{HEADS}_i - \text{TAILS}_i$
- $\text{HEADS}_i$ — $i$ 次投掷后头面向上的总次数
- $\text{TAILS}_i$ — $i$ 次投掷后背面向上的总次数

$$|X_i - X_{i-1}| \leq 1 \quad \Pr[|X_n| \geq tn^{1/2}] \leq 2 \exp(-t^2/2)$$

输赢次数之差的很大也是不太可能的!



## 5.2 鞅的一般形式

停时定理  
瓦尔德方程  
鞅尾不等式



鞅

**定义**

$Y_0, Y_1, Y_2, \dots$ 称为随机变量序列  $X_0, X_1, X_2, \dots$ 的鞅  
如果

$$Y_i \text{ 是 } X_0, X_1, X_2, \dots, X_i \text{ 的函数} \quad \forall i \geq 1$$

$$\mathbb{E}[Y_i | X_0, X_1, \dots, X_{i-1}] = Y_{i-1} \quad \forall i \geq 1$$



### 随机变量的和

设独立随机变量序列  $X_1, X_2, \dots$  满足  $E[X_i]=0$  ( $\forall i \geq 1$ )

定义

$$Y_i = X_1 + X_2 + \dots + X_i \quad \forall i \geq 1$$

$$E[Y_i | X_1, \dots, X_{i-1}]$$

$$= E[X_i + Y_{i-1} | X_1, \dots, X_{i-1}]$$

$$= E[X_i | X_1, \dots, X_{i-1}] + E[Y_{i-1} | X_1, \dots, X_{i-1}]$$

$$= E[X_i] + E[Y_{i-1} | X_1, \dots, X_{i-1}]$$

$$= 0 + Y_{i-1}$$

$$= Y_{i-1}$$

均值为0的随机变量之和是一个鞅

### 例1



### 随机变量和的平方

设  $X_1, X_2, \dots$  是均值为0, 方差为  $\sigma^2$  的独立随机变量

定义

$$Y_i = [X_1 + X_2 + \dots + X_i]^2 - i\sigma^2 \quad \forall i \geq 1$$

$$E[Y_i | X_1, \dots, X_{i-1}]$$

$$= E[X_i^2 + 2X_i(\sum_{k=1}^{i-1} X_k) + (\sum_{k=1}^{i-1} X_k)^2 - i\sigma^2 | X_1, \dots, X_{i-1}]$$

$$= E[X_i^2 - \sigma^2 | X_1, \dots, X_{i-1}] + 2E[X_i(\sum_{k=1}^{i-1} X_k) | X_1, \dots, X_{i-1}] + E[Y_{i-1} | X_1, \dots, X_{i-1}]$$

$$= (E[X_i])^2 + 2E[X_i]E[(\sum_{k=1}^{i-1} X_k) | X_1, \dots, X_{i-1}] + E[Y_{i-1} | X_1, \dots, X_{i-1}]$$

$$= Y_{i-1}$$

均值为0的随机变量和的平方是一个鞅

### 例2



### 随机变量函数的Doob序列

设  $f(X_1, X_2, \dots, X_n)$  是随机变量  $X_1, X_2, \dots, X_n$  的函数, 定义

$Y_i = E[f(X_1, X_2, \dots, X_n) | X_1, X_2, \dots, X_i] \quad \forall i=0, 1, \dots, n$

为函数  $f(X_1, X_2, \dots, X_n)$  的Doob序列

$$E[Y_i | X_1, \dots, X_{i-1}]$$

$$= E[E[f(X_1, X_2, \dots, X_n) | X_1, X_2, \dots, X_i] | X_1, \dots, X_{i-1}]$$

$$= E[E[f(X_1, X_2, \dots, X_n) | X_1, \dots, X_i] | X_1, \dots, X_{i-1}]$$

$$= Y_{i-1}$$

性质  $E[Y | Z] = E[E[Y | X, Z] | Z]$

Doob序列是一个鞅

### 杜比鞅



$$Y_0 = E[f(X_1, \dots, X_n)]$$

$$f(\underbrace{(\text{coin}, \text{coin}, \text{coin}, \text{coin}, \text{coin}, \text{coin})}_{\text{取均值}})$$

取均值



$$Y_1 = E[f(X_1, \dots, X_n) | X_1]$$

取定  $X_1$  的值

$$f(\underbrace{(1, \text{coin}, \text{coin}, \text{coin}, \text{coin}, \text{coin})}_{\text{取均值}})$$

取均值



$$Y_2 = E[f(X_1, \dots, X_n) | X_1, X_2]$$

取定  $X_1, X_2$  的值

$$f(\underbrace{(1, 0, \text{coin}, \text{coin}, \text{coin}, \text{coin})}_{\text{取均值}})$$

取均值



$$Y_3 = E[f(X_1, \dots, X_n) | X_1, X_2, X_3]$$

取定  $X_1, X_2, X_3$  的值

$$f(1, 0, 0, \underbrace{(\text{coin}), (\text{coin}), (\text{coin})}_{\text{取均值}})$$



$$Y_4 = E[f(X_1, \dots, X_n) | X_1, X_2, X_3, X_4]$$

取定  $X_1, X_2, X_3, X_4$  的值

$$f(1, 0, 0, 1, \underbrace{(\text{coin}), (\text{coin})}_{\text{取均值}})$$



$$Y_5 = E[f(X_1, \dots, X_n) | X_1, X_2, X_3, X_4, X_5]$$

取定  $X_1, X_2, X_3, X_4, X_5$  的值

$$f(1, 0, 0, 1, 0, \underbrace{(\text{coin})}_{\text{取均值}})$$



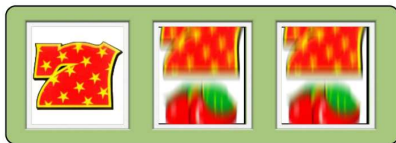
$$Y_n = E[f(X_1, \dots, X_n) | X_1, X_2, X_3, \dots, X_n]$$

取定  $X_1, X_2, X_3, \dots, X_n$  的值

$$f(1, 0, 0, 1, 0, 1)$$

$f(X_1, \dots, X_n)$  的一个具体取值

$$Y_n - Y_0 = f(X_1, \dots, X_n) - E[f(X_1, \dots, X_n)]$$



$Y_1$  = 看到第1个示数之后的平均收益



## 鞅的性质

**定理:** 如果  $Y_0, Y_1, Y_2, \dots$  是随机变量序列  $X_0, X_1, X_2, \dots$  的鞅, 则  $E[Y_i] = E[Y_0]$

**证明:** 因为  $Y_0, Y_1, Y_2, \dots$  是随机变量序列  $X_0, X_1, X_2, \dots$  的鞅,

$$Y_{i+1} = E[Y_i | Y_0, Y_1, \dots, Y_{i-1}]$$

两端同时求数学期望得,

$$E[Y_{i+1}] = E[E[Y_i | Y_0, Y_1, \dots, Y_{i-1}]] = E[Y_i]$$

$$E[Y_i] = E[Y_{i+1}] = \dots = E[Y_0]$$



## 鞅的停时定理

**定义：** 设  $Y_0, Y_1, Y_2, \dots$  是随机变量序列  $X_0, X_1, X_2, \dots$  的鞅，如果随机变量  $T=n$  仅依赖于  $Y_0, Y_1, Y_2, \dots, Y_n$  的取值，则称  $T$  是鞅  $\{Y_i | i \geq 0\}$  的一个**停时**。

**定理（鞅的停时定理）：** 设  $Y_0, Y_1, Y_2, \dots$  是随机变量序列  $X_0, X_1, X_2, \dots$  的鞅， $T$  是鞅  $\{Y_i | i \geq 0\}$  的一个停时，如果  $T$  是有限的，则  $E[Y_T] = E[Y_0]$ 。



## 应用1

### 公平赌局

- 初始赌资  $X_0 = b$ ， $X_i$  表示第  $i$  轮之后的赌资
- 每轮以  $1/2$  的概率赢，赢得的赌资是  $a$
- 以  $1/2$  的概率输，输掉的赌资是  $a$
- 每局的输赢是相互独立的
- 玩家赢得  $L_1$  或者输掉  $L_2$  之后停止游戏
- 问：玩家赢得  $L_1$  并停止游戏的概率  $q$  是多大？

游戏停止时间记为  $T$

$T=n \Leftrightarrow X_T - X_0 = L_1 \ (X_T \geq X_0)$  或  $X_T - X_0 = -L_2 \ (X_T < X_0)$  在第  $n$  轮首次发生

$T$  是鞅  $\{X_i\}$  的停时  $E[X_T] = E[X_0] = b$

$$(b+L_1)q + (b-L_2)(1-q) = b$$

$$q = L_2 / (L_1 + L_2)$$



## 应用2

### 选举定理

- 两人竞选， $A$  得  $a$  票， $B$  得  $b$  票 ( $a > b$ )
- 计票过程是  $a+b$  张选票的所有排列中均匀独立选取的
- 问：计票过程中  $A$  始终领先  $B$  的概率有多大？

$$n = a + b$$

$S_i$  = 统计  $i$  张选票之后  $A$  领先  $B$  的票数

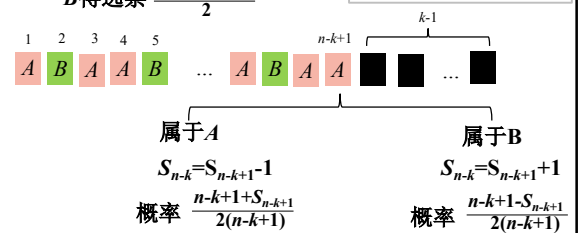
$$S_n = a - b$$



$$A \text{ 得选票 } \frac{n-k+1+S_{n-k+1}}{2}$$

$$B \text{ 得选票 } \frac{n-k+1-S_{n-k+1}}{2}$$

$$\begin{array}{ll} A \text{ 得 } x \text{ 票} & x+y = n-k+1 \\ B \text{ 得 } y \text{ 票} & x-y = S_{n-k+1} \end{array}$$



$$E[S_{n-k} | S_{n-k+1}] = S_{n-k+1} \frac{n-k}{n-k+1}$$

$$E\left[\frac{S_{n-k}}{n-k} \mid \frac{S_{n-k+1}}{n-k+1}\right] = \frac{S_{n-k+1}}{n-k+1}$$



### 选举定理

- 两人竞选， $A$  得  $a$  票， $B$  得  $b$  票 ( $a > b$ )
- 计票过程是  $a+b$  张选票的所有排列中均匀独立选取的
- 问：计票过程中  $A$  始终领先  $B$  的概率有多大？

$$n = a + b$$

$S_i$  = 统计  $i$  张选票之后  $A$  领先  $B$  的票数

$$S_n = a - b$$

$$E\left[\frac{S_{n-k}}{n-k} \mid \frac{S_{n-k+1}}{n-k+1}\right] = \frac{S_{n-k+1}}{n-k+1}$$

$$X_k = \frac{S_{n-k}}{n-k} \quad 0 \leq k \leq n-1$$

$$E[X_k | X_0, X_1, \dots, X_{k-1}] = X_{k-1}$$

$X_0, X_1, \dots, X_{n-1}$  是鞅



### 选举定理

- 两人竞选， $A$  得  $a$  票， $B$  得  $b$  票 ( $a > b$ )
- 计票过程是  $a+b$  张选票的所有排列中均匀独立选取的
- 问：计票过程中  $A$  始终领先  $B$  的概率有多大？

$$n = a + b$$

$S_i$  = 统计  $i$  张选票之后  $A$  领先  $B$  的票数

$$S_n = a - b$$

$$X_k = \frac{S_{n-k}}{n-k}$$

$X_0, X_1, \dots, X_{n-1}$  是鞅

$T=k$  是满足  $X_k=0$  的最小  $k$ ，如果存在这样的  $k$  值

$T=n-1$  如果不存在上述  $k$  值

$T=i$  仅依赖于  $X_0, \dots, X_i$  的取值，故  $E[X_T] = E[X_0] = (a-b)/(a+b)$



**第一种停时的特征:**  $T < n-1$

存在  $k$  使得  $X_k=0$ ,  $T$  就是这种  $k$  值的最小值。故  $X_T=0$

**第二种停时的特征:**  $T=n-1$

不存在  $k$  使得  $X_k=0$  且  $X_0=(a-b)/n > 0$ , 故  $X_T > 0$  恒成立

$$X_k = \frac{S_{n-k}}{n-k}$$

$S_1, S_2, \dots, > 0$  恒成立

A 得第1票并一直保持领先

$$X_{n-1} = X_T = S_1 = 1$$



**选举定理**

- 两人竞选, A 得  $a$  票, B 得  $b$  票 ( $a > b$ )
- 计票过程是  $a+b$  张选票的所有排列中均匀独立选取的
- 问: 计票过程中 A 始终领先 B 的概率有多大?  $\frac{a-b}{a+b}$

$$n = a+b$$

$S_i$  = 统计  $i$  张选票之后 A 领先 B 的票数  $S_n = a-b$

$$X_k = \frac{S_{n-k}}{n-k} \quad X_0, X_1, \dots, X_{n-1} \text{ 是鞅}$$

$T=k$  是满足  $X_k=0$  的最小  $k$ , 如果存在这样的  $k$  值

$T=n-1$  如果不存在上述  $k$  值

$T=i$  仅依赖于  $X_0, \dots, X_i$  的取值, 故  $E[X_T] = E[X_0] = (a-b)/(a+b)$

$$\Pr(\text{A 一直领先}) \cdot 1 + [1 - \Pr(\text{一直领先})] \cdot 0 = (a-b)/(a+b)$$



## 瓦尔德方程

**定理(瓦尔德方程):** 设  $X, X_1, X_2, \dots$  是独立同分布的随机变量,  $T$  是  $\{X_i | i \geq 1\}$  的一个停时, 如果  $E[T]$  和  $E[X]$  均是有限的, 则

$$E\left[\sum_{i=1}^T X_i\right] = E[T] \cdot E[X]$$

$\{X_i - E[X_i] \mid i \geq 1\}$  是均值为 0 的随机变量序列

$$Y_i = \sum_{j=1}^i (X_j - E[X_j]) \quad Y_1, Y_2, \dots \text{ 是鞅} \quad T \text{ 是停时}$$

$$E[Y_T] = E[Y_1] = 0 \quad \text{停时定理}$$

$$E[Y_T] = E\left[\sum_{i=1}^T X_i - T \cdot E[X]\right] = E\left[\sum_{i=1}^T X_i\right] - E[T] \cdot E[X] = 0$$



**两轮骰子赌局**

- 第一轮: 投掷均匀骰子的点数  $X$
- 第二轮: 投掷均匀骰子  $X$  次得点数  $Y_1, \dots, Y_X$
- 玩家收益  $Z = Y_1 + \dots + Y_X$
- 问: 玩家平均收益  $E[Z] = ?$



$$Y_1, Y_2, \dots \text{ 独立同分布} \quad E[Y_i] = 7/2$$

$$X \text{ 是随机变量序列 } Y_1, Y_2, \dots \text{ 的停时} \quad E[X] = 7/2$$

$$E[Z] = E\left[\sum_{i=1}^X Y_i\right] = E\left[\sum_{i=1}^X Y_i\right] - E[X] \cdot E[Y_i] = 0$$

$$E[Z] = E[X] \cdot E[Y_i] = 49/4$$



## 应用2: Las Vegas 算法的期望运行时间

**LAZYSELECT( $S, k$ )**

1.  $R$  = 独立、均匀、可放回地从  $S$  随机选取的  $n^{3/4}$  元素;
2. 在  $O(n)$  时间内排序  $R$ ;
3.  $x = (k/n)n^{3/4}$ ;
4.  $l = \max\{\quad, 0\}$ ;  $h = \min\{\quad, n^{3/4}\}$ ;
5.  $L = \min(R, l)$ ;  $H = \min(R, h)$ ;
6.  $L_p = \text{Rank}(S, L)$ ,  $H_p = \text{Rank}(S, H)$ ; (参见第2章)
7.  $P = \{y \in S \mid L \leq y \leq H\}$ ;
8. If  $\min(S, k) \in P$  and  $|P| \leq 4n^{3/4} + 1$
9. Then 排序  $P$ ,  $\min(S, k) = \min(P, (k - L_p))$ , 算法结束;
10. ELSE goto 1.

算法的期望运行时间是多少?



## Las Vegas 算法的期望运行时间

**Las Vegas 算法  $A(x)$**   
期望运行时间为  $T(n)$   
找到解的概率为  $p(n)$

**算法  $B(x)$**

1. while(true)
2.  $y = A(x)$
3. If  $y$  是问题的解 Then 返回  $y$

$Y_i$  是第  $i$  遍调用算法  $A$  的实际运行时间

$Y_1, Y_2, \dots$  是均值为  $T(n)$  的独立同分布的随机变量

算法终止时刻是其停时, 仅依赖于前面运行是否找到解

算法的期望运行时间 = 期望运行遍数  $\times T(n) = T(n)/p(n)$



HIT CS&E

### 应用3:共享总线服务器通信

#### 共享总线服务器通信

- $n$ 台服务器，共享总线，各自用缓冲区缓存消息
- 时间分片，每个时间片段各服务器至多发送一个消息
- 每个时间片
  - 各服务器以 $1/n$ 的概率发送一个消息
  - 有冲突则所有消息发送失败
  - 无冲突则消息发送成功
- 问：每台服务器都至少发送成功一个消息，需要多少时间片段？

在一个时间片段内，若只有一台服务器发消息，则成功

$$p = \binom{n}{1} \frac{1}{n} \left(1 - \frac{1}{n}\right)^{n-1} \approx e^{-1}$$

HIT CS&E

参数为 $p \approx e^{-1}$ 的几何分布  
每次成功前期望轮数为 $1/p = e$

参数为 $n, 1/n$ 的赠券收集过程  
 $n$ 台机器均发消息成功  
期望成功次数 =  $n \ln n$

所有机器均至少成功发送1次消息前的时间片段数  
= 期望成功次数  $\times$  每次成功前期望轮数  
=  $n \ln n \times e$   
=  $en \ln n$

请您构造随机变量序列，利用瓦尔德方程将分析过程精确化

HIT CS&E

### Azuma-Hoeffding不等式

#### Azuma不等式

如果 $Y_0, Y_1, Y_2, \dots$ 是随机变量序列 $X_0, X_1, X_2, \dots$ 的鞅，且对 $k \geq 1$ 满足 $|Y_k - Y_{k-1}| \leq c_k$ ，则

$$\Pr[|Y_n - Y_0| \geq t] \leq 2 \exp\left(-\frac{t^2}{2 \sum_{k=1}^n c_k^2}\right)$$

#### 推论：

如果 $Y_0, Y_1, Y_2, \dots$ 是随机变量序列 $X_0, X_1, X_2, \dots$ 的鞅且对 $k \geq 1$ 满足 $|Y_k - Y_{k-1}| \leq c$ ，则

$$\Pr[|Y_n - Y_0| \geq tcn^{1/2}] \leq 2 \exp\{-t^2/2\}$$

HIT CS&E

## 5.3 鞅的应用

### 5.3.1 模式匹配

### 5.3.2 球和箱子模型中空箱子个数

### 5.3.3 随机图的色数

HIT CS&E

### 5.3.1 模式匹配

HIT CS&E

### 模式匹配

#### 模式匹配问题

- 输入：大小为 $s$ 的字符集 $\Sigma$ 上长度为 $n$ 的串 $X=(X_1, \dots, X_n)$   
 $\Sigma$ 上长度为 $k$ 的模式串 $B=(B_1, \dots, B_k)$
- 输出： $B$ 在 $X$ 中的所有出现位置
- 问：这种匹配“有意义”吗？

大量的DNA片段看上去像随机片段  
不包含任何特殊的信息  
如果 $B$ 是这种随机片段，模式匹配将毫无意义

$B$ 有意义  $\Leftrightarrow$  它在随机选定的DNA序列的频率显著地偏离均值

那么一个固定的片段在随机DNA中的平均频率是多少呢？

HIT CS&E

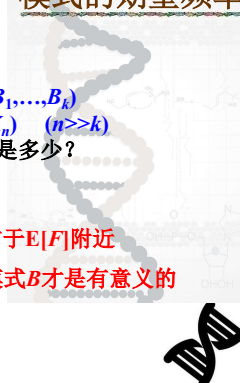
### 模式的期望频率

**模式的期望频率**

- 固定大小为 $s$ 的字符集 $\Sigma$
- 固定 $\Sigma$ 上长度为 $k$ 的模式串 $B=(B_1, \dots, B_k)$
- 随机取长度为 $n$ 的串 $X=(X_1, \dots, X_n)$  ( $n \gg k$ )
- 问:  $B$ 在 $X$ 出现频率 $F$ 的期望值是多少?

$$E[F] = (n-k+1)(1/s)^k$$

本节将证明:  $F$ 高概率地集中分布于 $E[F]$ 附近  
这意味着:  $F$ 频率偏离 $E[F]$ 后, 模式 $B$ 才是有意义的



HIT CS&E

### 模式的期望频率

$E[F] = (n-k+1)(1/s)^k$   $X_1, X_2, \dots, X_n, X_{i+1}$


↓ 在 $B$ 中的对齐位置

$Z_0 = E[F]$   
 $Z_i = E[F | X_1, X_2, \dots, X_i]$   $1 \leq i \leq n$   
 $Z_n = F$  即 $B$ 在 $X$ 中的出现次数

$|Z_{i+1} - Z_i| \leq k$  即  $|E[F | X_1, \dots, X_{i+1}] - E[F | X_1, \dots, X_i]| \leq k$

**推论:** 如果 $Y_0, Y_1, Y_2, \dots$ 是随机变量序列 $X_0, X_1, X_2, \dots$ 的鞅且对 $k \geq 1$ 满足  $|Y_k - Y_{k-1}| \leq c$   
 则  $\Pr[|Y_n - Y_0| \geq tcn^{1/2}] \leq 2 \exp\{-t^2/2\}$

$\Pr[|Z_n - Z_0| \geq \epsilon] \leq 2 \exp\{-\epsilon^2/(2nk^2)\}$



HIT CS&E

### 5.3.2 球和箱子模型中空箱子个数

HIT CS&E

### 重申球和箱子模型

$m$ 个球

均匀独立地将球投入箱子: 空箱子有几个?

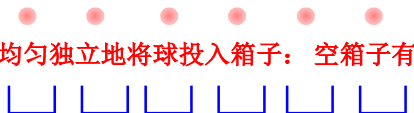
$n$ 个箱子

$X_i = \begin{cases} 1 & \text{第} i \text{个箱子是空的} \\ 0 & \text{否则} \end{cases}$   $\Pr[X_i=1] = (1-1/n)^m$

$X = \sum_{i=1}^n X_i$

$E[X] = n(1 - \frac{1}{n})^m$

下面证明:  $X$ 高概率地集中分布于 $E[X]$ 附近



HIT CS&E

$Y_i$ 表示第 $i$ 个球落入的箱子  $i=1, 2, \dots, n$   
 $X$ 表示所有球投掷完成后空箱子的个数

$Z_0 = E[X]$   
 $Z_i = E[X | Y_1, Y_2, \dots, Y_i]$   $1 \leq i \leq n$  **Doob鞅**  
 $Z_n = X$  即所有球投掷完成后空箱子的个数

前 $i$ 个球落入的箱子确定后

如果第 $i+1$ 个球落入空箱子, 则 $Z_{i+1} - Z_i = -1$   
 如果第 $i+1$ 个球未落入空箱子, 则 $Z_{i+1} - Z_i = 0$

$|Z_{i+1} - Z_i| \leq 1$   $\Pr[|Z_n - Z_0| \geq \epsilon] \leq 2 \exp\{-\epsilon^2/(2n)\}$

**推论:** 如果 $Y_0, Y_1, Y_2, \dots$ 是随机变量序列 $X_0, X_1, X_2, \dots$ 的鞅且对 $k \geq 1$ 满足  $|Y_k - Y_{k-1}| \leq c$   
 则  $\Pr[|Y_n - Y_0| \geq tcn^{1/2}] \leq 2 \exp\{-t^2/2\}$

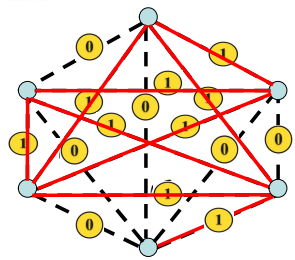
HIT CS&E

### 5.3.3 随机图的色数

**随机图模型**

$G(n, p)$

- $n$  个顶点
- 任意顶点对成边的概率为  $p$



**随机图的参数**

$G(n, p)$

随机图的参数  $f(G)$

- 直径: 最短路径的最大值
- 围长: 最短环的边数
- 色数: 真着色所需最少颜色数

Doob 鞅  
边暴露鞅

估计  $f(G)$  集中性的一般方法

将所有可能的边依次编号为  $1, 2, \dots, \binom{n}{2}$ ,  $X_i$  表示第  $i$  边的取舍

$$Y_0 = E[f(G)]$$

$$Y_i = E[f(G) | X_1, \dots, X_i]$$

$$Y_{n(n-1)/2} = E[f(G) | X_1, \dots, X_p, \dots, X_{n(n-1)/2}] = f(G)$$

$$\Pr[|f(G) - E[f(G)]| > t] = \Pr[|Y_{n(n-1)/2} - Y_0| > t] = ?$$

**随机图的色数**

$G(n, p)$

随机图的色数  $f(G)$

- 色数: 真着色所需最少颜色数
- $\chi(G)$

将所有可能的边依次编号为  $1, 2, \dots, \binom{n}{2}$ ,  $X_i$  表示第  $i$  边的取舍

$$X_i = \begin{cases} 1 & \text{第 } i \text{ 条边在 } G \text{ 中} \\ 0 & \text{否则} \end{cases}$$

$$Y_0 = E[\chi(G)]$$

$$Y_i = E[\chi(G) | X_1, \dots, X_i]$$

$$Y_{n(n-1)/2} = \chi(G)$$

**边暴露鞅示例**

$Y_i = E[\chi(G) | X_1, \dots, X_i]$

$Y_0 = E[\chi(G)]$

$Y_1$

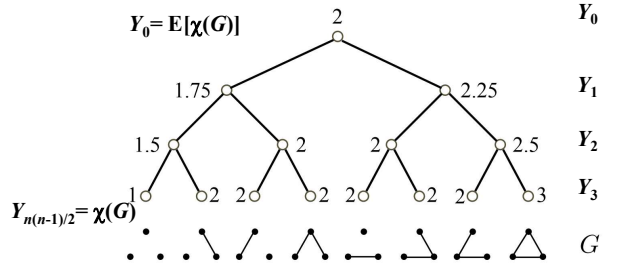
$Y_2$

$Y_3$

$Y_{n(n-1)/2} = \chi(G)$

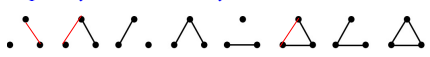
$G$

容易直接验证:  $E[Y_{i+1} | X_1, \dots, X_i] = Y_i$



**随机图的色数**

$X_1, \dots, X_i$  取定之后得  $G_i$



在  $G_i$  上根据  $X_{i+1}$  添加第  $i+1$  条边

情形1:  $X_{i+1}=0$ ,  $\chi(G_{i+1}) = \chi(G_i)$

情形2:  $X_{i+1}=1$ , 第  $i+1$  条边连接两个无边顶点  
 $\chi(G_{i+1}) \leq \chi(G_i) + 1$

情形3:  $X_{i+1}=1$ , 第  $i+1$  条边连接无边顶点和有边顶点  
 $\chi(G_{i+1}) \leq \chi(G_i) + 1$

情形4:  $X_{i+1}=1$ , 第  $i+1$  条边连接两个有边顶点  
 $\chi(G_{i+1}) \leq \chi(G_i) + 1$

$|Y_{i+1} - Y_i| \leq 1$

**随机图的色数**

$Y_0 = E[\chi(G)]$

$Y_i = E[\chi(G) | X_1, \dots, X_i]$

$Y_{n(n-1)/2} = \chi(G)$

$|Y_{i+1} - Y_i| \leq 1$

**推论:** 如果  $Y_0, Y_1, Y_2, \dots$  是随机变量序列  $X_0, X_1, X_2, \dots$  的鞅且对  $k \geq 1$  满足  $|Y_k - Y_{k-1}| \leq c$   
则  $\Pr[|Y_n - Y_0| \geq tcn^{1/2}] \leq 2 \exp\{-t^2/2\}$

$\Pr[|\chi(G) - E[\chi(G)]| \geq t(n(n-1)/2)^{1/2}] \leq 2 \exp\{-t^2/2\}$