

# Число различных элементов

Всеволод Опарин

CS Клуб, Осень 2014

19 Октября

## Базовые понятия

- ▶ Случайная величина (С.В.)  $X$ .
- ▶  $\mathbf{E}[X] = \sum_{i \in \text{Dom}(X)} i \cdot \mathbf{Pr}[X = i]$ .
- ▶ Для набора С.В.  $X_1, X_2, \dots, X_n$  и чисел  $\alpha_1, \dots, \alpha_n$   
 $\mathbf{E}[\sum_i \alpha_i \cdot X_i] = \sum_i \alpha_i \cdot \mathbf{E}[X_i]$
- ▶  $\mathbf{D}[X] = \mathbf{E}[(X - \mathbf{E}[X])^2] = \mathbf{E}[X^2] - \mathbf{E}[X]^2$ .
- ▶  $\mathbf{D}[\alpha X] = \alpha^2 \mathbf{D}[X]$ .

## Базовые понятия

- ▶ С.В.  $X$  и  $Y$  независимы, если  $\forall a, b \in \mathbb{R}$

$$\Pr[X = a \wedge Y = b] = \Pr[X = a] \cdot \Pr[Y = b].$$

- ▶ Для независимых С.В.  $X$  и  $Y$

$$\mathbf{E}[XY] = \mathbf{E}[X]\mathbf{E}[Y],$$

$$\mathbf{D}[X + Y] = \mathbf{D}[X] + \mathbf{D}[Y].$$

- ▶ Набор С.В.  $X_1, \dots, X_n$  2-независим, если для любых  $1 \leq i \neq j \leq n$   $X_i$  и  $X_j$  независимы.

- ▶  $X_1, X_2, X_3 \in \{0, 1\}$ ,  $X_1 + X_2 + X_3 = 0 \pmod{2}$ .

- ▶ Для 2-независимого набора С.В.

$$\mathbf{D}[\sum_i X_i] = \sum_i \mathbf{D}[X_i].$$

## Основные неравенства

- ▶  $\Pr[A \vee B] \leq \Pr[A] + \Pr[B]$ .
- ▶  $\Pr[A \wedge B] \leq \Pr[A]$ .

### Неравенство Маркова

Для С.В.  $X \geq 0$  и числа  $a > 0$

$$\Pr[X \geq a] \leq \frac{\mathbf{E}[X]}{a}.$$

# Основные неравенства

## Неравенство Чебышёва

Для С.В.  $X$  при  $\mathbf{E}[X] = \mu$

$$\Pr[|X - \mu| \geq \varepsilon] \leq \frac{\mathbf{D}[X]}{\varepsilon^2}.$$

## Неравенство Чернова

Для набора независимых С.В.  $X_1, X_2, \dots, X_n$  таких, что  $X_i \in \{0, 1\}$  и  $\mathbf{E}[X_i] = \mu$ , и константы  $\varepsilon > 0$

$$\Pr\left[\left|\frac{1}{n} \cdot \sum_i X_i - \mu\right| \geq \varepsilon\right] \leq \exp(-c \cdot n).$$

## Хэш-функции

- ▶ Пусть  $\mathcal{H} = \{h_i\}_{i=1}^t$ , где  $h_i : K \rightarrow V$ .  $\mathcal{H}$  – семейство 2-независимых хэш-функций, если для любых  $k_1 \neq k_2$  и любых  $v_1, v_2$

$$\Pr_{h \leftarrow U(\mathcal{H})} [ h(k_1) = v_1 \wedge h(k_2) = v_2 ] = \frac{1}{|V|^2}.$$

- ▶ Пример  $h_i : \mathbb{F}_p \rightarrow \mathbb{F}_p$ .  $h_{a,b}(x) = a \cdot x + b$ .

$$\Pr_{h \leftarrow U(\mathcal{H})} [ a \cdot k_1 + b = v_1 \wedge a \cdot k_2 + b = v_2 ] = \frac{1}{p^2}.$$

$$\begin{cases} a \cdot k_1 + b = v_1, \\ a \cdot k_2 + b = v_2. \end{cases}$$

# Число различных элементов

## Задача

- ▶  $\sigma = \langle a_1, \dots, a_m \rangle$ ,  $a_i \in [n] = [2^l]$ .
- ▶  $f[x] = \#\{i \mid a_i = x\}$  – частота.
- ▶  $d = \#\{x \mid f[x] > 0\}$  – число различных элементов.
- ▶ Разрешен один проход. Найти  $d$ .
  
- ▶ FM'83 AMS'99:  $\Pr [ans \in [\frac{d}{3}, 3d]] \geq 1 - \delta$ .  
Память  $O(\log \frac{1}{\delta} \log n)$ .
- ▶ VJKST'04:  $\Pr [ |ans - d| \leq \epsilon d ] \geq 1 - \delta$   
Память  $O(\log \frac{1}{\delta} (\log n + \frac{1}{\epsilon^2} (\log \frac{1}{\epsilon} + \log \log n)))$ .

- ▶  $\text{zero}(p) = \max\{i \mid p \text{ делится на } 2^i\}$ .
- ▶ Алгоритм:
  1.  $\text{init}()$ :  $z = 0$ , взять  $h : [n] \rightarrow [n]$  из 2-независимого семейства.
  2.  $\text{process}(y) : z = \max(z, \text{zero}(h(y)))$ .
  3.  $\text{answer}()$ : вернуть  $2^{z+\frac{1}{2}}$ .
- ▶ Пусть  $d$  – ответ.  $\bar{d}$  – результат алгоритма.  
 $\Pr [\bar{d} \geq 3 \cdot d] \leq 0.47$ .  
 $\Pr [\bar{d} \leq \frac{1}{3}d] \leq 0.47$ .



## AMS. Анализ

▶  $X_{j,r} = [\text{zero}(h(j)) \geq r].$

$$\mathbf{E}[X_{j,r}] = \frac{1}{2^r},$$

$$\mathbf{D}[X_{j,r}] \leq \frac{1}{2^r}.$$

▶  $Y_r = \sum_{j|f[j]>0} X_{j,r}.$

$$\mathbf{E}[Y_r] = \frac{d}{2^r},$$

$$\mathbf{D}[Y_r] \leq \frac{d}{2^r}.$$

▶  $\Pr[Y_r > 0] = \Pr[Y_r \geq 1] \leq \mathbf{E}[Y_r]/1 = \frac{d}{2^r}.$

▶  $\Pr[Y_r = 0] \leq \Pr[|Y_r - \frac{d}{2^r}| \geq \frac{d}{2^r}] \leq \frac{\mathbf{D}[Y_r]}{(d/2^r)^2} \leq \frac{2^r}{d}.$

- ▶ Пусть  $a$  – минимальное такое, что  $3 \cdot d \leq 2^{a+\frac{1}{2}}$ . Тогда

$$\Pr[ Y_a > 0 ] \leq \frac{d}{2^a} \leq \frac{\sqrt{2}}{3}$$

- ▶ Пусть  $b$  – максимальное такое, что  $2^{b+\frac{1}{2}} \leq \frac{1}{3} \cdot d$ . Тогда

$$\Pr[ Y_{b+1} = 0 ] \leq \frac{2^{b+1}}{d} \leq \frac{\sqrt{2}}{3}$$

# Медиана

- ▶ Алгоритм возвращает значение, большее разрешенного, с вероятностью  $p < \frac{1}{2}$ .
- ▶ Запустим алгоритм независимо  $k$  раз и возьмем медиану.
- ▶  $T_i =$  [ошибка при  $i$ -ом запуске].
- ▶  $\mathbf{E}[T_i] \leq p$ .
- ▶ Если медиана плохая, то алгоритм ошибся хотя бы половину раз.
- ▶  $\Pr \left[ \left| \frac{1}{k} \sum_i T_i - \mathbf{E}[T_i] \right| \geq \frac{1}{2} - p \right] \leq \exp(-c \cdot k)$ .
- ▶ Запустим  $\log(\frac{2}{\delta})$  и получим ошибку  $\frac{\delta}{2}$  с одной стороны.
- ▶  $\Pr$  [оценка мала или велика]  $\leq$   
 $\Pr$  [мала] +  $\Pr$  [велика]  $\leq \delta$

## VJKST. Алгоритм

```
init():  
    B = []  
    z = 0  
    pick 2-ind h : [n] -> [n]  
  
process(y):  
    zy = zero(h(y))  
    if zy >= z:  
        B <- (y, zy)  
        while (B > c / (eps * eps)):  
            z++  
            remove all (a, b) from B s.t. b < z  
  
ans = |B| * pow(2, z)
```

## ВКСТ. Анализ

- ▶  $|B| = Y_z$ .
- ▶  $\mathbf{E}[|B|] = \mathbf{E}[Y_z] = \frac{d}{2^z}$ .
- ▶  $\mathbf{D}[|B|] = \mathbf{D}[Y_z] \leq \frac{d}{2^z}$ .
- ▶ При  $z = 0$  всё точно.
- ▶  $\left| |B| \cdot 2^z - d \right| \geq \varepsilon \cdot d \Leftrightarrow \left| |B| - \frac{d}{2^z} \right| \geq \varepsilon \cdot \frac{d}{2^z}$

## ВКСТ. Анализ

- ▶ Пусть  $s$  такое, что  $\frac{12}{\epsilon^2} \leq \frac{d}{2^s} \leq \frac{24}{\epsilon^2}$ .
- ▶  $\Pr \left[ \left| |B| - \frac{d}{2^z} \right| \geq \epsilon \cdot \frac{d}{2^z} \right] =$   
 $\sum_{i=1}^{\log n} \Pr \left[ \left| |B| - \frac{d}{2^z} \right| \geq \epsilon \cdot \frac{d}{2^z} \wedge z = i \right] \leq$   
 $\sum_{i=1}^{s-1} \Pr \left[ \left| |B| - \frac{d}{2^z} \right| \geq \epsilon \cdot \frac{d}{2^z} \right] + \sum_{i=s}^{\log n} \Pr [z = i] \leq \frac{1}{6}$
- ▶  $\Pr \left[ \left| |B| - \frac{d}{2^i} \right| \geq \epsilon \cdot \frac{d}{2^i} \right] \leq \frac{\mathbf{D}[|B|]}{(\epsilon d / 2^i)^2} \leq \frac{d}{2^i} \frac{1}{\epsilon^2} \left( \frac{2^i}{d} \right)^2 = \frac{1}{\epsilon^2} \frac{2^i}{d}$ .
- ▶  $\sum_{i=1}^{s-1} \Pr \left[ \left| |B| - \frac{d}{2^z} \right| \geq \epsilon \cdot \frac{d}{2^z} \right] \leq \frac{1}{\epsilon^2} \frac{2^s}{d} \leq \frac{1}{12}$ .
- ▶  $\sum_{i=s}^{\log n} \Pr [z = i] = \Pr \left[ Y_{s-1} \geq \frac{c}{\epsilon^2} \right] \leq \frac{d}{2^{s-1}} \cdot \frac{\epsilon^2}{c} \leq \frac{24 \cdot 2}{\epsilon^2} \cdot \frac{\epsilon^2}{c} \leq$   
 $\frac{48}{c}$ .

## ВКСТ. Память

- ▶  $O(\frac{1}{\epsilon^2} \log \frac{1}{\delta} \log n)$ .
- ▶ Было:  
 $B \leftarrow (y, \text{zero}(h(y)))$ .
- ▶ Стало:
  1. В *init*: вытаскиваем еще одну 2-независимую функцию  $g : [n] \rightarrow [b \cdot \frac{1}{\epsilon^4} \log^2 n]$ .
  2. В *process*:  $B \leftarrow (g(y), \text{zero}(h(y)))$ .
  3.  $\Pr[\text{коллизия}] \leq \frac{1}{6}$ .
- ▶ Память:  $O(\log \frac{1}{\delta} (\log n + \frac{1}{\epsilon^2} (\log \frac{1}{\epsilon} + \log \log n)))$

Спасибо!

Вопросы?