



哈爾濱工業大學(深圳)
HARBIN INSTITUTE OF TECHNOLOGY, SHENZHEN

随 机 算 法

第2章 矩与离差

2.1 理论回顾

定理 2.1 期望的线性性

对于任意一组有限个具有有限期望的离散型随机变量 X_1, X_2, \dots, X_n , 有

$$E\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n E[X_i]$$

引理 2.2

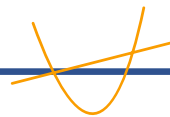
对于任意的常数 c 和离散型随机变量 X , 有

$$E[cX] = cE[X]$$

2.1 理论回顾

凸函数定义: $\forall x_1, x_2, x_1 \neq x_2$

$\forall \lambda, 0 \leq \lambda \leq 1$



定理 2.3 詹森不等式

如果 f 是一个凸函数, 那么

凸函数引理: 若 f 二次可微, f 凸 $\Leftrightarrow f'' \geq 0$

$$f(\lambda x_1 + (1-\lambda)x_2) \leq \lambda f(x_1) + (1-\lambda)f(x_2)$$

$$E[f(X)] \geq f(E[X])$$

注释:

更一般的例子: $E[X^2] \geq (E[X])^2$ 。

证明:
$$f(x) = f(x_0) + f'(x_0)(x-x_0) + \frac{f''(\xi)}{2} (x-x_0)^2$$

$$f(x) \geq f(x_0) + f'(x_0)(x-x_0)$$

$$E[f(x)] \geq E[f(x_0)] + f'(x_0)E[(x-x_0)]$$

令 $x_0 = E[x]$, 则

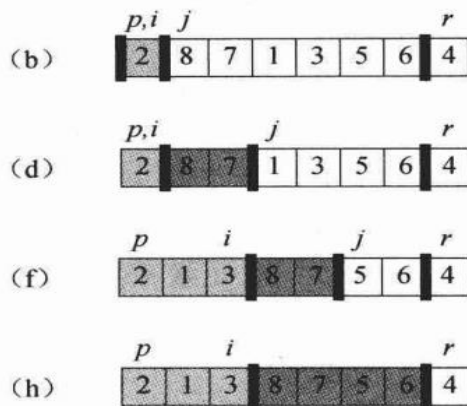
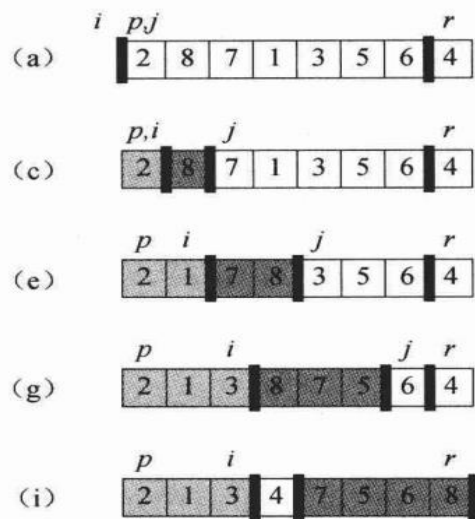
$$E[f(x)] \geq \underbrace{E[f(E[x])]}_{\downarrow \text{常数}} + \underbrace{f'(x_0)[E[x] - E[x_0]]}_{\text{等于0}}$$

$$E[f(x)] \geq f(E[X])$$

2.1 理论回顾

■ 应用：快速排序的期望运行时间

输入是 n 个不同的数 x_1, x_2, \dots, x_n 的列表，选择一个基准元素 x ，将其他每个元素与 x 进行比较，进而分成两个子表：小于 x 的元素和大于 x 的元素。快速排序法是对这些子表的递归排序。



2.1 理论回顾

算法 2.1 快速排序

输入：全序总体上 n 个不同元素的列表 $S = \{x_1, x_2, \dots, x_n\}$.

输出：排序后的 S 的元素.

1. 如果 S 只有一个或零个元素，返回 S ；否则继续.
2. 选择 S 中一个元素作为基准元素，称为 x .
3. 为了将其他元素分成两个子列表， S 中的每个其他元素与 x 作比较；
 - a) S_1 是 S 中所有比 x 小的元素
 - b) S_2 是 S 中所有比 x 大的元素
4. 对 S_1 和 S_2 进行快速排序.
5. 返回列表 S_1, x, S_2 .

最差比较次数 $(n-1) + (n-2) + \dots + 1 = \frac{n(n-1)}{2}$

2.1 理论回顾

- 把算法改成随机选取基准，将快速排序变成一种随机化算法。

① 随机选取基准 x_i

② 逐一比较，分为子集 S_1, S_2

③ 对 S_1, S_2 进行排序

期望次数为 $2n \ln n + O(n)$ ，期望是关于基准的随机选取。

- 可能性是保持原有的确定性算法，用列表的第一个元素作为基准，但考虑输入的概率模型。

期望次数为 $2n \ln n + O(n)$ ，期望是关于输入的概率模型。



2.1 理论回顾

定理 2.4

假设在随机快速排序法中，每一次都是从所有可能中独立且随机地选取基准的，那么对于任意的输入，随机快速排序法所做比较的期望次数为 $2n \ln n + O(n)$

n 个数，期望比较次数 $E(n)$

$n=0, 1$ 时, $E(n)=0$, $n \geq 2$ 时

$$E(n) = \frac{1}{n-1} \sum_{k=1}^{n-1} [E(k) + E(n-k) + n]$$

$$= \frac{2}{n-1} \sum_{k=1}^{n-1} E(k) + n$$

$$\text{记 } S(n) = \sum_{k=1}^n E(k)$$

$$S(n) - S(n-1) = \frac{2}{n-1} S(n-1) + n$$

$$S(n) = \frac{n+1}{n-1} S(n-1) + n$$

$$\frac{1}{n(n+1)} S(n) = \frac{1}{n(n-1)} S(n-1) + \frac{1}{n+1}$$

$$\text{记 } t(n) = \frac{1}{n(n+1)} S(n)$$

$$t(n) = t(n-1) + \frac{1}{n+1}$$

$$t(1) = 0;$$

$$t(2) = \frac{1}{3};$$

$$t(3) = \frac{1}{3} + \frac{1}{4};$$

...

$$t(n-1) = \sum_{i=1}^{n-1} \frac{1}{i} - \frac{1}{2} - 1 \\ = H(n) - \frac{3}{2}$$

其中 $H(n)$ 为调和级数
 $\ln n \leq H(n) \leq \ln n + 1$
 $H(n) = \ln n + O(1)$

$$t(n) = \frac{1}{n+1} + H(n) - \frac{3}{2}$$

$$= H(n+1) - \frac{3}{2}$$

\Downarrow

$$S(n) = n(n+1)t(n)$$

$$E(n) = S(n) - S(n-1)$$

$$= n(n+1)t(n) - n(n-1)t(n-1)$$

$$\Downarrow t(n) = t(n-1) + \frac{1}{n+1}$$

$$= 2nt(n-1) + n$$

$$= 2nH(n) - 2n$$

$$= 2n \ln n + \underbrace{2nO(1)}_{O(n)} - 2n$$

$$= 2n \ln n + O(n) \quad \checkmark$$

定理 2.5

假设在随机快速排序法中，每次选取子列表中第一个元素作为基准，如果输入是在其所有可能排列中均匀随机选取的，那么确定性快速排序法所做比较的期望次数为 $2n \ln n + O(n)$ 。

2.1 理论回顾

定理 2.6 马尔可夫不等式

设 X 是只取非负值的随机变量, 那么对所有 $a > 0$, 有

$$\Pr(X \geq a) \leq \frac{E[X]}{a}$$

构造 $I(x) = \mathbb{1}(x \geq a)$, 一定有 $I(x) \leq \frac{x}{a}$

$$E(I) = P(x \geq a) \leq E\left(\frac{x}{a}\right) = \frac{E(x)}{a}$$

注释:

重要性在于不知道 X 的分布 $(f(x), p_k)$ 情况下, 通过 $E[X]$ 估计事件 $\{X \geq \varepsilon\}$ 的概率上限。

2.1 理论回顾

定义 2.1

- ◆ 一个随机变量 X 的 k 阶矩为 $E[X^k]$.
- ◆ 一个随机变量 X 的方差和标准差分别定义为

$$\text{Var}[X] = E[(X - E[X])^2] = E[X^2] - (E[X])^2 \text{ 和 } \overset{\text{标准差}}{\sigma[X]} = \sqrt{\text{Var}[X]}$$

- ◆ 两个随机变量 X 和 Y 的协方差为
- $$\text{Cov}(X, Y) = E[(X - E[X])(Y - E[Y])]$$

定理 2.7

对任意两个随机变量 X 和 Y ，有

$$\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y] + 2\text{Cov}(X, Y)$$

2.1 理论回顾

定理 2.8 切比雪夫不等式

对任意的 $a > 0$, 有

$$\Pr(|X - E[X]| \geq a) \leq \frac{\text{Var}[X]}{a^2}$$

马尔科夫不等式

$|X - E[X]|$ 也是一个 r.v.
 $\Pr(|X - E[X]| \geq a)$
 $= \Pr((X - E[X])^2 \geq a^2)$
 $\leq \frac{E[(X - E[X])^2]}{a^2} = \frac{\text{Var}[X]}{a^2}$

推论 2.9 直接由切比雪夫不等式可得

对任意的 $t > 1$, 有 $\Pr(|X - E[X]| \geq t \cdot \sigma[X]) \leq \frac{1}{t^2}$

$$\Pr(|X - E[X]| \geq t \cdot E[X]) \leq \frac{\text{Var}[X]}{t^2 (E[X])^2}$$

2.1 理论回顾

■ 问题描述

假设有 n 种不同的赠券，每盒麦片内附有其中的一种赠券，每种赠券获取机率相同，而且无限供应。

- 假定每盒麦片中的赠券是从 n 中可能中独立且随机选取的
- 假定收集赠券时不与其他人合作

在拥有每种赠券至少一张之前，需要购买多少盒麦片？



2.1 理论回顾

■ 令 X 表示收集到每种赠券至少一张所需要购买的麦片盒数，求收集 n 张赠券的时间 X 的期望为 $E[X]$ 。

■ 如果 X_i 表示恰有 $i - 1$ 种不同的赠券时所购买的盒数，那么得到 $X = \sum_{i=1}^n X_i$ ，

当得到恰好 $i - 1$ 种赠券时，再得到一张新赠券的概率是 $p_i = 1 - \frac{i-1}{n}$ ，因此

$$E[X_i] = \frac{1}{p_i} = \frac{n}{n-i+1} \longrightarrow E[X] = E[\sum_{i=1}^n X_i] = \sum_{i=1}^n E[X_i] = n \sum_{i=1}^n \frac{1}{i}$$

X_i 服从几何分布 $Ge(p_i)$, $p_i = \frac{n-i+1}{n}$

$$D(X_i) = \frac{1-p_i}{p_i^2} = \frac{n(i-1)}{(n-i+1)^2}$$

$$D(X) = \sum_{i=1}^n D(X_i) = \sum_{i=1}^n \frac{n(i-1)}{(n-i+1)^2}$$

$$\leq \sum_{i=1}^n \frac{n^2}{(n-i+1)^2} = n^2 \sum_{i=1}^n \frac{1}{i^2}$$

$$\leq \frac{\pi^2}{6} n^2 \leftarrow \sum_{i=1}^{\infty} \frac{1}{i^2} = \frac{\pi^2}{6}$$

称为调和数 H_n

2.1 理论回顾

引理 2.10 调和数 $H_n = \sum_{i=1}^n \frac{1}{i}$ 满足 $H_n = \ln n + O(1)$.

由马尔可夫不等式得 $\Pr(X \geq 2nH_n) \leq \frac{E[X]}{2nH_n} = \frac{1}{2}$

引理 2.11 参数为 p 的几何随机变量的方差是 $\frac{1-p}{p^2}$. 期望为 $\frac{1}{p}$

利用切比雪夫不等式, 得到 切比雪夫不等式可以给出更高的精度

$$\Pr(X \geq 2nH_n) \Leftrightarrow \Pr(|X - nH_n| \geq nH_n) \leq \frac{\text{Var}[X]}{(nH_n)^2} = \frac{\pi^2}{6(H_n)^2} = O\left(\frac{1}{\ln^2 n}\right)$$

$$\ln n \leq H_n \leq \ln n + 1$$


2.2 中位数和平均值

定义 2.2 中位数 随机变量的中位数

设 X 是随机变量, X 的中位数定义为满足下列条件的 m 值

$$\Pr(X \leq m) \geq 1/2 \text{ 和 } \Pr(X \geq m) \geq 1/2$$

$x_1, x_2, \dots, x_{2k+1}$  x_{k+1} 中位数是 x_{k+1}

x_1, x_2, \dots, x_{2k}  (x_k, x_{k+1}) 中位数是区间 (x_k, x_{k+1}) 任意数

2.2 中位数和平均值

定理 2.12 对于具有有限期望 $E[X]$ 和有限中位数 m 的任意随机变量 X

1. 期望 $E[X]$ 是使下列表达式最小的 c 的值 $E[(X - c)^2]$
2. 中位数 m 是使下列表达式最小的 c 的值 $E[|X - c|]$

证明: 1. $E[(X - c)^2]$ 2. 记住就行

$$= E[X^2 - 2cX + c^2]$$

$$= E[X^2] - 2cE[X] + c^2$$

开口向上抛物线, 当且仅
当 $c = E[X]$ 时取最小

2.2 中位数和平均值

定理 2.13 如果 X 是随机变量，且具有有限标准差 σ ，期望 μ 和中位数 m ，
则 $|\mu - m| \leq \sigma$ (证明期末考)

证明： $|\mu - m|$
 $= |\mathbb{E}(X) - m|$
 $= |\mathbb{E}(X - m)|$
 $\leq \mathbb{E}(|X - m|)$ \searrow 詹森不等式，选取 $|x|$ 为凸函数
 $\leq \mathbb{E}(|X - \mu|)$ \searrow 中位数性质
 $\leq \sqrt{\mathbb{E}[(X - \mu)^2]} = \sigma$ \searrow 詹森不等式，选取 $\sqrt{|x|^2}$ 为凸函数

2.3 应用

■ 计算中位数的随机化算法

目的：找到两个元素 d 和 u ，它们依 S 的排序是彼此接近的，且中位数 m 于它们之间。

即：1. $d \leq m \leq u$

→ C 的规模要比较小

2. 对 $C = \{s \in S: d \leq s \leq u\}$, $|C| = O(n/\log n)$ (在 d 和 u 之间的元素总数是少的)

主要思想：从 S 中有放回的抽样，得到一个含有 $\lceil n^{3/4} \rceil$ 个元素的多重集合 R 。以大概率保证集合 C

包含中位数 m ，固定 d 和 u 分别为 R 的排序的第 $\lfloor \frac{1}{2}n^{3/4} - \sqrt{n} \rfloor$ 个和第 $\left(\lceil \frac{1}{2}n^{3/4} + \sqrt{n} \rceil\right)$ 个元素。保证

a. 集合 C 足够大，以大概率包含 m ； b. 集合 C 足够小，以大概率用次线性时间排序。

算法 2.2 随机化中位数算法

输入：一个全序总体上 n 个元素的集合 S .

输出： S 的中位数元素，用 m 表示.

1. 独立地、均匀随机地，有放回地从 S 中取出 $\lceil n^{3/4} \rceil$ 个元素组成一个（多重）集合 R .
2. 对集合 R 排序.
3. 设 d 为排序集合 R 中第 $\left(\left\lfloor \frac{1}{2} n^{3/4} - \sqrt{n} \right\rfloor\right)$ 个最小元素.
4. 设 u 为排序集合 R 中第 $\left(\left\lceil \frac{1}{2} n^{3/4} + \sqrt{n} \right\rceil\right)$ 个最小元素.
5. 将集合 S 中每个元素与 d 和 u 比较，计算集合 $C = \{x \in S: d \leq x \leq u\}$ 及数 $\ell_d = |\{x \in S: x < d\}|$ 和 $\ell_u = |\{x \in S: x > u\}|$.
6. 如果 $\ell_d > n/2$ 或 $\ell_u > n/2$,则输出FAIL.
7. 如果 $|C| \leq 4n^{3/4}$ ，则对集合 C 排序；否则，输出FAIL. 对 C 进行控制,保障运行时间不过长
8. 输出排序集合 C 中的第 $(\lfloor n/2 \rfloor - \ell_d + 1)$ 个元素.

定理 2.14 随机化中位数算法以线性时间结束，而且如果它输出的不是 FAIL，则它输出的是输入集合 S 的正确中位数元素

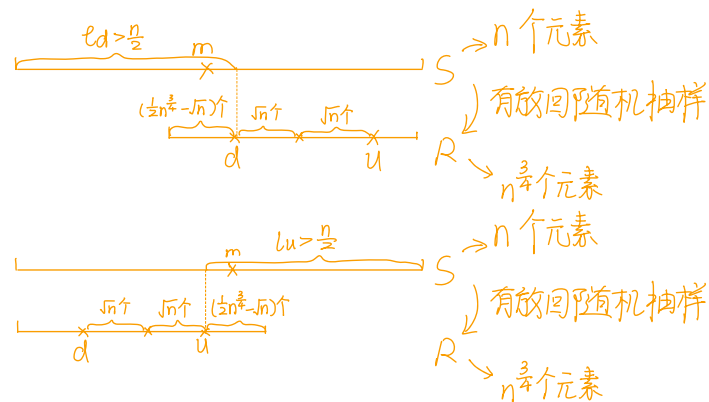
定理 2.15 随机化中位数算法失败的概率的界为 $n^{-1/4}$

$$IP\{Fail\} = IP\{l_d > \frac{n}{2}\} + IP\{l_u > \frac{n}{2}\} + IP\{|c| > 4n^{\frac{3}{4}}\} \leq n^{-\frac{1}{4}}$$

- 中位数的蒙特卡罗随机化算法可以通过重复运行直到成功为止而转化为 **Las Vegas** 算法，也就是说把它转化为 **Las Vegas** 算法意味着运行时间是可变的，尽管期望运行时间是线性的。

中位数算法失败

- $l_d > \frac{n}{2} \Leftrightarrow d > m \Leftrightarrow R$ 中元素 \leq 中位数 m 的个数小于 $\frac{1}{2}n^{\frac{3}{4}} - \sqrt{n}$
- $l_u > \frac{n}{2} \Leftrightarrow u < m \Leftrightarrow R$ 中元素 \geq 中位数 m 的个数小于 $\frac{1}{2}n^{\frac{3}{4}} - \sqrt{n}$
- $|c| > 4n^{\frac{3}{4}}$



■ 考虑下面三个事件

$$\varepsilon_1: Y_1 = |\{r \in R | r \leq m\}| < \frac{1}{2}n^{3/4} - \sqrt{n} \Leftrightarrow \text{算法在第6步的 } \ell_d > \frac{n}{2} \text{ 失败}$$

$$\varepsilon_2: Y_2 = |\{r \in R | r \geq m\}| < \frac{1}{2}n^{3/4} - \sqrt{n} \Leftrightarrow \text{算法在第6步的 } \ell_u > \frac{n}{2} \text{ 失败}$$

$$\varepsilon_3: |C| > 4n^{3/4} \Leftrightarrow \text{算法在第七步失败}$$

$$\Pr(\varepsilon_1) \leq \frac{1}{4}n^{-1/4}$$

$$\Pr(\varepsilon_2) \leq \frac{1}{4}n^{-1/4}$$

$$\Pr(\varepsilon_3) \leq \frac{1}{2}n^{-1/4}$$

证明在下一页

$$\begin{aligned} \Pr\{\text{Fail}\} &= \Pr\{\ell_d > \frac{n}{2}\} + \Pr\{\ell_u > \frac{n}{2}\} + \Pr\{|C| > 4n^{3/4}\} \\ &= \Pr\{\varepsilon_1\} + \Pr\{\varepsilon_2\} + \Pr\{\varepsilon_3\} \leq n^{-1/4} \end{aligned}$$

证明:

$$(1) \mathbb{P}(\mathcal{E}_1) \leq \frac{1}{4}n^{-\frac{1}{4}}$$

定义随机变量 $X_i = \begin{cases} 1, & \text{第 } i \text{ 次抽样} \leq \text{中位数 } m \\ 0, & \text{其它} \end{cases}$

$$P\{X_i = 1\} = \frac{\frac{(n-1)}{2} + 1}{n} = \frac{1}{2} + \frac{1}{2n}, \text{ 事件 } \varepsilon_1 \Leftrightarrow Y_1 = \sum_{i=1}^n \frac{1}{2} < \frac{1}{2}n^{\frac{3}{2}} - \sqrt{n}$$

Y_1 是伯努利试验的和, Y_1 服从二项分布, $Y_1 \sim B(n^{\frac{3}{4}}, \frac{1}{2} + \frac{1}{2n})$

$$\text{Var}(Y_i) = n^{\frac{3}{4}} \left(\frac{1}{2} + \frac{1}{2n} \right) \left(\frac{1}{2} - \frac{1}{2n} \right) = \frac{1}{4} n^{\frac{3}{4}} - \frac{1}{4} n^{-\frac{5}{4}} < \frac{1}{4} n^{\frac{3}{4}}$$

由切比雪夫不等式:

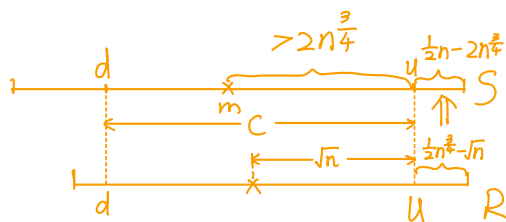
$$P(\mathcal{E}_1) = P\left\{Y_1 < \frac{1}{2}n^{\frac{2}{3}} - \sqrt{n}\right\} = P\{|Y_1 - E(Y_1)| > \sqrt{n}\} \leq \frac{\text{Var}[Y_1]}{n} < \frac{1}{4}n^{-\frac{1}{4}}$$

$$\textcircled{2} \quad IP(\varepsilon_2) \leq \frac{1}{4} n^{-\frac{1}{4}}$$

证明同①

$$\textcircled{3} \mathbb{P}(\mathcal{E}_3) \leq \frac{1}{2} n^{-\frac{1}{2}}$$

若 ε_3 发生, 则以下两事件至少有一个发生:

$$E_{3.1} = C \text{ 中至少 } 2n^{\frac{3}{4}} \text{ 个元素大于中位数}$$
$$E_{3,2} = C \text{ 中至少 } 2n^{\frac{2}{3}} \text{ 个元素大于中位数}$$


考虑 $\varepsilon_{3.1}$ 发生的概率的界:

若 C 中至少 $2n^{\frac{3}{4}}$ 个元素大于中位数, 依 S 从小到大的排序, u 的次序 $\geq \frac{1}{2}n + 2n^{\frac{3}{4}}$

这样 R 中至少有 $\frac{1}{2}n - \sqrt{n}$ 次抽样是在 S 的 $\frac{1}{2}n - 2n^{\frac{3}{4}}$ 个最大元素中进行的

定义随机变量 $X_i = \begin{cases} 1, & \text{第 } i \text{ 次抽样是在 } S \text{ 的 } \frac{1}{2}n - 2n^{\frac{2}{3}} \text{ 个最大元素中} \\ 0, & \text{其它} \end{cases}$

$$P\{X_i = 1\} = \frac{1}{2} - 2n^{-\frac{1}{4}}$$

$$\text{令 } X = \sum_{i=1}^{n^{\frac{3}{4}}} X_i \sim B(n^{\frac{3}{4}}, \frac{1}{2} - 2n^{-\frac{1}{4}})$$

$$E[X] = \frac{1}{2}n^{\frac{3}{4}} - 2\sqrt{n}$$

$$\text{Var}[X] = n^{\frac{3}{4}} \left(\frac{1}{2} - 2n^{-\frac{1}{4}} \right) \left(\frac{1}{2} + 2n^{-\frac{1}{4}} \right) = \frac{1}{4} n^{\frac{3}{4}} - 4n^{\frac{1}{4}} < \frac{1}{4} n^{\frac{3}{4}}$$

由切比雪夫不等式:

$$P\{\varepsilon_{n,1}\} = P\{X \geq \frac{1}{2}n^{\frac{2}{3}} - \sqrt{n}\} \leq P\{|X - E[X]| \geq \sqrt{n}\} \leq \frac{\text{Var}[X]}{n} < \frac{1}{4}n^{-\frac{1}{2}}$$

同理: $|\mathcal{P}\{\varepsilon_{3,2}\}| < \frac{1}{4}n^{-\frac{1}{4}}$

从而: $|\mathcal{P}\{\mathcal{E}_3\}| \leq |\mathcal{P}\{\mathcal{E}_{3.1}\}| + |\mathcal{P}\{\mathcal{E}_{3.2}\}| < \frac{1}{2}n^{-\frac{1}{4}}$