

大数据算法

Big Data Algorithms

▣ 概率基础

▣ 亚线性空间算法

▣ 亚线性时间算法

▣ 并行模型算法

3 / 250

概率基础

定义1.1[概率空间, probability space]

一个概率空间包含三要素:

- ▷ 样本空间(sample space) Ω : 概率空间描述的随机行为所有可能结果
- ▷ 事件(events) \mathcal{F} : 每个事件是 Ω 的一个子集合
- ▷ 概率函数 $\Pr: \mathcal{F} \rightarrow [0, 1]$: 满足概率测度要求

满足概率测度需要如下条件

- ▷ 事件集合 \mathcal{F} 包含 Ω , 对集合补操作和集合并操作封闭

- ▷ 概率函数满足

- 非负性(\Pr 取值非负)
- 规范性($\Pr(\emptyset) = 0, \Pr(\Omega) = 1$)
- 完全可加性(对任意两两不相交的事件 A_1, \dots, A_n 有 $\Pr(\cup A_i) = \sum_i \Pr(A_i)$), 事件序列是无穷可数的也成立

在通常的离散概率空间(discrete probability space), 每一个 Ω 中的元素都是一个基本事件, $\mathcal{F} = 2^\Omega$

4 / 250

概率基础

定义1.2[条件概率, conditional probability]

给定事件 E_2 , E_1 发生的条件概率为

$$\Pr(E_1|E_2) = \frac{\Pr(E_1 \cap E_2)}{\Pr(E_2)}$$

条件概率只有当 $\Pr(E_2) > 0$ 的时候才有意义

定义1.3[全概率公式, law of total probability]

令 E_1, E_2, \dots, E_n 为概率空间的两两不相交事件且 $\cup E_i = \Omega$, 我们有

$$\Pr(B) = \sum_{i=1}^n \Pr(B \cap E_i) = \sum_{i=1}^n \Pr(B|E_i)\Pr(E_i)$$

5 / 250

概率基础

定义1.6[随机变量, random variable]

一个样本空间 Ω 上的随机变量 X 是一个实值函数 $X: \Omega \mapsto \mathcal{R}$, 一个离散随机变量是取值为有穷或者可数无穷多个值的变量

- ▷ 给定离散随机变量 X 和实数 a , 事件 $X = a$ 代表集合 $\{s \in \Omega: X(s) = a\}$

$$\Pr(X = a) = \sum_{s \in \Omega: X(s)=a} \Pr(s)$$

定义1.7[离散随机变量期望]

$$\mathbb{E}[X] = \sum_i i \cdot \Pr(X = i)$$

定理1.2[期望的另一种计算方式]

X 是一个取值为非负整数的离散随机变量, 有 $\mathbb{E}[X] = \sum_{i=1}^{\infty} \Pr(X \geq i)$

12 / 250

概率基础

定义1.8[连续随机变量, continuous random variable]

一个随机变量 X 的**概率分布**由它的**分布函数** $F(x)$ 给出, 对任意 $x \in \mathbb{R}$, 定义 $F(x) = \Pr(X \leq x)$; X 是连续随机变量, 如果 $F(x)$ 是 x 的连续函数。

- ▷ 连续随机变量 X , 必有 $\Pr[X = x] = 0$, 因此有 $\Pr(X \leq x) = \Pr(X < x)$

定义1.9[概率密度函数]

如果存在函数 $f(x)$, 有

$$F(a) = \int_{-\infty}^a f(t) dt$$

, 则称 $f(x)$ 是 $F(x)$ 的密度函数, 且 $f(x) = F'(x)$

定理1.3[连续随机变量的期望]

X 一个连续随机变量, 有 $\mathbb{E}[X] = \int_{-\infty}^{+\infty} xf(x) dx$

13 / 250

概率基础

定理1.4[马尔可夫不等式]

对任意非负随机变量 X 和 $a > 0$, 有 $\Pr(X \geq a) \leq \frac{\mathbb{E}[X]}{a}$

定义1.11[方差]

$$\text{Var}[X] = \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$$

定理1.5[切比雪夫不等式]

对任意随机变量 X 和 $a > 0$, 有 $\Pr(|X - \mathbb{E}[X]| \geq a) \leq \frac{\text{Var}[X]}{a^2}$

定理1.6[切尔诺夫不等式, Chernoff/Hoeffding Bound]

令 X_1, X_2, \dots, X_m 是独立的、取值 $\in \{0, 1\}$ 的随机变量, 变量和的期望为 $\mu = \mathbb{E}[\sum_i X_i]$, 令 $\epsilon \in [0, 1]$, 有 $\Pr[|\sum_i X_i - \mu| > \epsilon \mu] \leq 2e^{-\epsilon^2 \mu / 3}$.

16 / 250

▣ 概率基础

▣ 亚线性空间算法

▣ 亚线性时间算法

▣ 并行模型算法

28 / 250

问题1：近似计数

流模型 (Streaming Model)

A stream:

$$\sigma = \langle a_1, a_2, \dots, a_n \rangle, a_i \in [m],$$

then we can define a frequency vector

$$\mathbf{f} = (f_1, \dots, f_m), \text{ s.t. } \sum_{1 \leq i \leq m} f_i = n.$$

定义1.13[The Counting Problem]

给定数据流 σ 以及 a_i , 计算 f_i 。

29 / 250

问题1：近似计数

为了计算 f_i , 假设 a_i 共出现 n 次

- ▷ 需要 $O(\log n)$ 位的空间代价存储 n
- ▷ 如果 n 非常大, $O(\log n)$ 也许不可接受
- ▷ Can we maintain the count in less space?
当然不可能, 除非我们放宽要求
- ▷ 我们将设计一个占用 $O(\log \log n)$ 位空间的近似算法

定义1.14[The Approximate Counting Problem]

给定数据流 σ 以及 a_i , 计算 \hat{f}_i , 满足

$$\Pr[|\hat{f}_i - f_i| > \epsilon f_i] < \delta$$

实际应用中, 可令 $\epsilon = 1/3$, $\delta = 1\%$

30 / 250

问题1：近似计数

Morris Algorithm

*/** Maintain a counter X **/*

1. $X \leftarrow 0$;
2. If the event a_i happens,
Update $X \leftarrow X + 1$ with prob. $\frac{1}{2^X}$
3. return $\hat{f}_i = 2^X - 1$

Analysis of Morris Algorithm

- ▷ 计算期望 $E[2^X - 1] = n$
- ▷ 计算方差 $\text{var}[2^X - 1] \leq \frac{3n(n+1)}{2} + 1 = O(n^2)$
- ▷ 利用切比雪夫不等式

32 / 250

问题1：近似计数

定理1.10[切比雪夫不等式, Chebyshev Inequality]

对任意的 $\lambda > 0$, $\Pr[|X - E[X]| > \lambda] \leq \frac{\text{var}[X]}{\lambda^2}$.

定理1.11[推论, Chebyshev Inequality Corollary]

$X = E[X] \pm \sqrt{\text{var}[X]/c}$ 成立的概率至少为 $1 - c$.

定理1.12[Morris算法分析]

至少以 $1 - c$ 的概率, Morris算法满足

$$2^{X_n} - 1 = E[2^{X_n} - 1] \pm \sqrt{\text{var}[2^{X_n} - 1]/c} = n \pm O(n)$$

35 / 250

问题1：近似计数

如何获得一个 $1 + \epsilon$ 近似算法

Morris+ Algorithm

*/** 重复 Morris 算法 k 次 **/*

1. Run k independent copies of Morris's algorithm.
2. Let the results be (X_1, \dots, X_k) .
3. return $\frac{1}{k} \sum_{i=1}^k (2^{X_i} - 1)$.

Analysis of Morris+ Algorithm

- ▷ 计算期望 $E[Y] = E[\frac{1}{k} \sum_{i=1}^k (2^{X_i} - 1)] = n$
- ▷ 计算方差 $\text{var}[Y] = O(\frac{n^2}{k})$
- ▷ 利用切比雪夫不等式

36 / 250

问题1：近似计数

定理1.13[Analysis of Morris+ Algorithm]

至少以 $1 - c$ 的概率, Morris+算法满足

$$Y = n \pm \frac{1}{\sqrt{k}} O(n)$$

- ▷ 令 $k = O(\frac{1}{\epsilon^2})$, 上述算法是一个 $(1 \pm \epsilon)$ 近似算法.
 - ▷ 以常数概率如 $\frac{9}{10}$ 满足近似标准.
- 进一步, 可以将概率改进到 $1 - \delta$, 对任意 $\delta > 0$.

38 / 250

问题1：近似计数

- ▷ 为了以 $1 - \delta$ 概率得到 $(1 \pm \epsilon)$ -近似, Morris+需要 $k = O(\frac{1}{\epsilon^2} \frac{1}{\delta})$ 次运行
- ▷ 假设, 有 Morris+算法以常数概率 $c \geq 0.9$ 得到 $1 \pm \epsilon$ 近似
- ▷ 设计一个 $1 - \delta$ 概率 $(1 \pm \epsilon)$ -近似算法, 代价为 $O(\frac{1}{\epsilon^2} \log \frac{1}{\delta} \log \log n)$

Morris++ Algorithm

*/** Median Trick **/*

1. 重复 Morris+算法 $t = O(\log(1/\delta))$ 次
2. 取 t 次估计结果的中间值 (median) 返回

Analysis of Morris++ Algorithm

- ▷ 计算期望
- ▷ 利用切尔诺夫不等式

40 / 250

问题1：近似计数

定理1.14[切尔诺夫不等式, Chernoff/Hoeffding Bound]

令 X_1, X_2, \dots, X_m 是独立的、取值在 $\{0, 1\}$ 的随机变量, 变量和的期望为 $\mu = E[\sum_i X_i]$, 令 $\epsilon \in [0, 1]$, 有 $\Pr[|\sum_i X_i - \mu| > \epsilon \mu] \leq 2e^{-\epsilon^2 \mu/3}$.

令随机变量 $X_i = 1 \Leftrightarrow i$ 个 Morris+算法运行的结果在正确范围内

$$E[X_i] = 0.9 \text{ 且 } \mu = 0.9t$$

易知, 如果 $\sum_i X_i > 0.5t$, 那么 Morris++ 算法将返回正确范围值

$$\begin{aligned} \Pr[\sum_i X_i < 0.5t] &= \Pr[\sum_i X_i - 0.9t < -0.4t] \\ &\leq \Pr[|\sum_i X_i - \mu| > 0.4t] \\ &= \Pr[|\sum_i X_i - \mu| > \mu 0.4/0.9] \end{aligned}$$

$$\begin{aligned} \text{利用切尔诺夫不等式} \Rightarrow &\leq 2e^{-\frac{16}{81 \times 3} \cdot 0.9t} \\ &\leq e^{-c't} \end{aligned}$$

若想要 $\Pr[\sum_i X_i < 0.5t] \leq e^{-c't} < \delta$, 只需要 $t = O(\log(1/\delta))$

41 / 250

问题2：不重复元素数

流模型 (Streaming Model)

A stream:

$$\sigma = \langle a_1, a_2, \dots, a_n \rangle, a_i \in [m],$$

then we can define a frequency vector

$$\mathbf{f} = (f_1, \dots, f_m), \text{ s.t. } \sum_{1 \leq i \leq m} f_i = n.$$

定义1.15 [The Distinct Elements Problem]

给定数据流 σ , 计算 $\sum_{1 \leq i \leq m} \mathbf{I}[f_i > 0]$, $\mathbf{I}[f_i > 0] = 1$ 当且仅当 $f_i > 0$ 。

42 / 250

问题2：不重复元素数

精确计算需要多少空间代价?

- ▷ 方法1: $O(m)$ 位.
为每一个 $[m]$ 中的元素维护一个位, 长度为 m 的向量.
- ▷ 方法2: $O(n \log(m))$ 位.
维护 n 个数, 每一个使用 $\log(m)$ 位.

43 / 250

问题2：不重复元素数

- ▷ 一个理想化的解决方案: 假设可以存储实数
- ▷ 利用哈希函数 (hash function)
 - $h: [m] \mapsto [0, 1]$
 - $h(i)$ 的函数值是 $[0, 1]$ 实数, 均匀分布

FM Algorithm [Flajolet-Martin 1985]

*/** Maintain a variable z **/*

1. 随机选取一个哈希函数 $h: [m] \mapsto [0, 1]$
2. $z = 1$.
3. 每遇到一个元素 i , 更新: $z = \min(z, h(i))$
4. return $1/z - 1$.

53 / 250

问题2：不重复元素数

FM算法的分析: 令 $X = \frac{1}{z} - 1$, $E[z] = \frac{1}{D+1}$, $\text{var}[z] \leq \frac{2}{(D+1)(D+2)}$

$$\Pr[|z - \frac{1}{D+1}| > c \frac{1}{D+1}] < \frac{2(D+1)^2}{c^2 \cdot (D+1)(D+2)} < \frac{2}{c^2}$$

$$\Pr[|X - D| > \epsilon D]$$

$$\leq \Pr[|z - \frac{1}{D+1}| > \frac{\epsilon D}{(D+1 + \epsilon D)} \frac{1}{(D+1)}]$$

$$\begin{aligned} c &= \frac{\epsilon D}{D+1 + \epsilon D} \\ &< \frac{2(D+1 + \epsilon D)^2}{\epsilon^2 D^2} \\ &= 2\left(\frac{D+1}{\epsilon D} + 1\right)^2 \\ &< 2\left(\frac{2}{\epsilon} + 1\right)^2 \end{aligned}$$

58 / 250

问题2：不重复元素数

$$\begin{aligned} &\Pr[|X - D| > \epsilon D] \\ \Leftrightarrow &\Pr[|1/z - 1 - D| > \epsilon D] \\ \Leftrightarrow &\Pr[\frac{1}{z} < 1 + (1 - \epsilon)D \text{ 或 } \frac{1}{z} > 1 + (1 + \epsilon)D] \\ \Leftrightarrow &\Pr[z > \frac{1}{1 + (1 - \epsilon)D} \text{ 或 } z < \frac{1}{1 + (1 + \epsilon)D}] \\ \Leftrightarrow &\Pr[z - \frac{1}{1 + D} > \frac{\epsilon D}{(D+1)(1 + D - \epsilon D)} \\ &\text{或 } z - \frac{1}{1 + D} < \frac{-\epsilon D}{(D+1)(1 + D + \epsilon D)}] \\ \leq &\Pr[|z - \frac{1}{D+1}| > \frac{\epsilon D}{1 + D + \epsilon D} \frac{1}{D+1}] \end{aligned}$$

59 / 250

问题2：不重复元素数

利用多次运行

FM+ Algorithm

*/** Maintain a variable z **/*

1. for j from 1 to k
2. 随机选取一个哈希函数 $h_j: [m] \mapsto [0, 1]$
3. $z_j = 1$.
4. 每次遇到 i , 更新: $z_j = \min(z_j, h_j(i))$
5. $Z = \frac{1}{k} \sum_{j=1}^k z_j$
6. return $1/Z - 1$.

60 / 250

问题2：不重复元素数

FM+算法分析: $X = \frac{1}{z} - 1$, $Z = \frac{1}{k} \sum_{j=1}^k z_j$, $E[Z] = E[z]$, $\text{var}[Z] = \frac{\text{var}[z]}{k}$

$$\Pr[|Z - \frac{1}{D+1}| > c \frac{1}{D+1}] < \frac{(D+1)^2}{c^2} \text{var}[Z]$$

$$c = \frac{\epsilon D}{D+1 + \epsilon D} \Rightarrow \Pr[|X - D| > \epsilon D] < \frac{2(D+1 + \epsilon D)^2}{k \epsilon^2 D^2} < \frac{2}{k} \left(\frac{2}{\epsilon} + 1\right)^2$$

这里, ϵ 为精度要求, 假设概率要求为 $1 - \delta$, 只需

$$\begin{aligned} &\frac{2}{k} \left(\frac{2}{\epsilon} + 1\right)^2 < \delta \\ \Rightarrow &k > \frac{2}{\delta} \left(\frac{2}{\epsilon} + 1\right)^2 = O\left(\frac{1}{\epsilon^2 \delta}\right) \end{aligned}$$

61 / 250

问题2：不重复元素数

利用Median技术

FM++ Algorithm

1. 运行具有常数概率的FM+算法 $m = \Theta(\log \frac{1}{\delta})$ 次;
2. \hat{D} 为所有 m 个结果的中间值;
3. return \hat{D} .

FM++ 是一个以 $1 - \delta$ 概率保证 $(1 \pm \epsilon)$ 近似的算法
整体代价为 $O(\frac{1}{\epsilon^2} \log \frac{1}{\delta})$

63 / 250

问题2：不重复元素数

利用多次运行，另一种做法

FM'+ (Bottom-k) Algorithm

*/** Maintain a variable z **/*

1. 随机选取一个哈希函数 $h: [m] \mapsto [0, 1]$
2. $(z_1, z_2, \dots, z_k) = 1$.
3. 维护当前看到的最小的 k 个哈希值，
满足 $z_1 \leq z_2 \leq \dots \leq z_k$
4. **return** k/z_k .

FM'+是一个 $(1 \pm \epsilon)$ 近似算法

$k = O(\frac{1}{\epsilon^2}) \Rightarrow$ FM'+以常数概率成功

65 / 250

问题2：不重复元素数

FM'+算法的分析

$$\Pr[|\frac{k}{z_k} - D| > \epsilon D] < ?$$

- ▷ $|\frac{k}{z_k} - D| > \epsilon D \Leftrightarrow z_k < \frac{k}{(1+\epsilon)D}$ 或者 $z_k > \frac{k}{(1-\epsilon)D}$
- ▷ 令随机变量 $X_i = 1 \Leftrightarrow h(i) < \frac{k}{(1+\epsilon)D}$
 - $z_k < \frac{k}{(1+\epsilon)D} \equiv \sum_{i=1}^m X_i > k$
- ▷ 令随机变量 $Y_i = 1 \Leftrightarrow h(i) < \frac{k}{(1-\epsilon)D}$
 - $z_k > \frac{k}{(1-\epsilon)D} \equiv \sum_{i=1}^m Y_i < k$
- ▷ $\mathbb{E}[\sum X_i] = \frac{k}{(1+\epsilon)} \Leftarrow m = D$
- ▷ $\text{var}[\sum X_i] = D \cdot \text{var}[X_i] \leq D \cdot \mathbb{E}[X_i^2] = D \cdot \mathbb{E}[X_i] = \frac{k}{1+\epsilon}$
- ▷ $\mathbb{E}[\sum Y_i] = \frac{k}{(1-\epsilon)}$
- ▷ $\text{var}[\sum Y_i] = D \cdot \text{var}[Y_i] \leq D \cdot \mathbb{E}[Y_i^2] = D \cdot \mathbb{E}[Y_i] = \frac{k}{1-\epsilon}$

66 / 250

问题2：不重复元素数

FM'+算法的分析

$$\begin{aligned} \Pr[\sum X_i > k] &= \Pr[\sum X_i - \frac{k}{1+\epsilon} > k - \frac{k}{1+\epsilon}] \\ &= \Pr[\sum X_i - \frac{k}{1+\epsilon} > \frac{\epsilon k}{1+\epsilon}] \\ &\leq \Pr[|\sum X_i - \frac{k}{1+\epsilon}| > \frac{\epsilon k}{1+\epsilon}] \\ &\leq \frac{(1+\epsilon)^2 \text{var}[\sum X_i]}{\epsilon^2 k^2} \leq \frac{1+\epsilon}{\epsilon^2 k} \leq \frac{c}{2} \end{aligned}$$

$$\Rightarrow k \geq \frac{2(1+\epsilon)}{\epsilon^2}$$

$$\begin{aligned} \Pr[\sum Y_i < k] &= \Pr[\sum Y_i - \frac{k}{1-\epsilon} < k - \frac{k}{1-\epsilon}] \\ &= \Pr[\sum Y_i - \frac{k}{1-\epsilon} < -\frac{\epsilon k}{1-\epsilon}] \\ &\leq \Pr[|\sum Y_i - \frac{k}{1-\epsilon}| > \frac{\epsilon k}{1-\epsilon}] \\ &\leq \frac{(1-\epsilon)^2 \text{var}[\sum Y_i]}{\epsilon^2 k^2} \leq \frac{1-\epsilon}{\epsilon^2 k} \leq \frac{c}{2} \end{aligned}$$

$$\Rightarrow k \geq \frac{2(1-\epsilon)}{\epsilon^2}$$

$$k \geq \max\{\frac{2(1-\epsilon)}{\epsilon^2}, \frac{2(1+\epsilon)}{\epsilon^2}\} = O(\frac{1}{\epsilon^2})$$

67 / 250

问题2：不重复元素数

事实上，我们无法存储实数

- ▷ 设计一个 $O(\log m)$ 位空间代价的 $O(1)$ 近似算法
- ▷ $\frac{1}{C} \times D \leq \hat{D} \leq C \times D$

Practical FM Algorithm

1. 从 2-wise independent 哈希函数族中随机选取 $h: [m] \mapsto [m]$;
2. $z = 0$;
3. **if** $\text{zeros}(h(j)) > z$ **then** $z = \text{zeros}(h(j))$;
 $\text{zeros}(p) = \max\{i \mid p \% 2^i = 0\}$ 01101010 $\Rightarrow \text{zeros}(\cdot) = 2$
4. **return** $\hat{D} = 2^{z+1/2}$;

73 / 250

问题2：不重复元素数

定义 1.16 [k-wise independent hash functions]

一个从 $[a]$ 到 $[b]$ 的哈希函数族 \mathcal{H} 被称作是 k -wise independent 的，如果 $\forall j_1, \dots, j_k \in [b]$ 和 $\forall i_1, \dots, i_k \in [a]$ ，其中 $i_x \neq i_y$ ，有

$$\Pr_{h \in \mathcal{H}}(h(i_1) = j_1 \wedge \dots \wedge h(i_k) = j_k) = 1/b^k$$

- ▷ 理论上，我们可以利用 $\log_2 |\mathcal{H}|$ 位存储某个 $h \in \mathcal{H}$

例 1.2 [k-wise independent 例子]

- ▷ 所有形如 $[a] \rightarrow [b]$ 的函数构成的 \mathcal{H} ， $|\mathcal{H}| = b^a$;
- ▷ $\mathcal{H}_{\text{poly}}$: 至多为 $k-1$ 阶多项式，系数在 $[n]$ 中，要求 $n = p^r$ ， p 为素数， $|\mathcal{H}_{\text{poly}}| = n^k$ ，存储需 $O(k \log n)$ 位
- ▷ 存储 2-wise 的 $\mathcal{H}_{\text{poly}}$ 某个函数，需要 $O(2 \log n)$ 位

74 / 250

问题2：不重复元素数

存储哈希函数需要 $O(\log m)$ 位；存储 z 需要 $\log \log m$ 位

定理 1.22

Practical FM 算法以 $1 - 2\sqrt{2}$ 的概率满足 $D/C \leq \hat{D} \leq C \times D$

- ▷ 对 $j \in [m]$ 和 $r \geq 0$ ，令 $X_{r,j} = 1 \Leftrightarrow \text{zeros}(h(j)) \geq r$
- ▷ $Y_r = \sum_{j \in [m]} X_{r,j} = \sum_{j: f_j > 0} X_{r,j}$
 - Y_r : 哈希函数值的二进制形式尾部 0 数目 $\geq r$ 的元素数
- ▷ 假设算法结束时， $z = t$
 - $t \geq r \Leftrightarrow$ 存在 j 使得 $\text{zeros}(h(j)) \geq r \Leftrightarrow X_{r,j} = 1 \Leftrightarrow Y_r > 0$
 - $t \leq r-1 \Leftrightarrow \forall j$ 有 $\text{zeros}(h(j)) \leq r-1 \Leftrightarrow X_{r,j} = 0 \Leftrightarrow Y_r = 0$

$$\mathbb{E}[X_{r,j}] = \Pr(\text{zeros}(h(j)) \geq r) = \frac{1}{2^r}$$

$$\mathbb{E}[Y_r] = \sum_{j \in [m]} \mathbb{E}[X_{r,j}] = \frac{D}{2^r}$$

75 / 250

问题2：不重复元素数

$$\begin{aligned} \text{var}[Y_r] &= \text{var}[\sum X_{r,j}] \\ &= \mathbb{E}[(\sum X_{r,j})^2] - \mathbb{E}[\sum X_{r,j}]^2 \\ &= \sum_{i,j} \mathbb{E}[X_{r,i} \cdot X_{r,j}] - \sum_{i,j} (\mathbb{E}[X_{r,i}] \cdot \mathbb{E}[X_{r,j}]) \\ \text{2-wise-indep.} &\equiv i \neq j \Rightarrow \mathbb{E}[X_{r,i} \cdot X_{r,j}] = \mathbb{E}[X_{r,i}] \cdot \mathbb{E}[X_{r,j}] \\ &= \sum \mathbb{E}[X_{r,j}^2] - \sum \mathbb{E}[X_{r,j}]^2 \\ &\leq \sum \mathbb{E}[X_{r,j}^2] = \sum \mathbb{E}[X_{r,j}] = \frac{D}{2^r} \end{aligned}$$

由 Markov 不等式有

$$\Pr[Y_r > 0] = \Pr[Y_r \geq 1] \leq \frac{\mathbb{E}[Y_r]}{1} = \frac{D}{2^r}$$

由切比雪夫不等式有

$$\Pr[Y_r = 0] \leq \Pr[|Y_r - \mathbb{E}[Y_r]| \geq \frac{D}{2^r}] \leq \frac{\text{var}[Y_r]}{(D/2^r)^2} \leq \frac{2^r}{D}$$

79 / 250

问题2：不重复元素数

再一次，利用 Median 技术

- ▷ 运行有常数概率 Practical FM 算法 $t = \Theta(\log \frac{1}{\delta})$ 次
- ▷ \hat{D} 为所有 t 个结果的中间值
- ▷ 上述算法以 $1 - \delta$ 概率保证 $O(1)$ 近似的算法
- ▷ 整体代价为 $O(\log m \log \frac{1}{\delta})$

81 / 250

问题2：不重复元素数

- ▷ $O(\log m + \frac{1}{\epsilon^2}(\log \frac{1}{\epsilon} + \log \log m))$ 空间代价的算法
- ▷ 主要思想：利用两层hash

BJKST Algorithm

1. 随机选2-wise independent 哈希函数 $h: [m] \rightarrow [m]$;
2. 随机选2-wise independent 哈希函数 $g: [m] \rightarrow [b\epsilon^{-4} \log^2 m]$;
3. $z = 0, B = \emptyset$;
4. **if** $\text{zeros}(h(j)) \geq z$ **then**
5. $B = B \cup \{(g(j), \text{zeros}(h(j)))\}$;
6. **if** $|B| > \frac{\epsilon}{\epsilon^2}$ **then**
7. $z = z + 1$;
8. 从 B 中删除 (α, β) , 其中 $\beta < z$;
9. **return** $\hat{D} = |B|2^z$;

82 / 250

问题2：不重复元素数

定理1.23[BJKST算法性能]

BJKST算法使用 $O(\log m + \frac{1}{\epsilon^2}(\log \frac{1}{\epsilon} + \log \log m))$ 空间，可以实现至少 $2/3$ 概率保证 $(1 \pm \epsilon)$ 近似。

- ▷ 存储 h 需 $O(\log m)$
- ▷ 存储 $|B|$ 个 g 取值需 $O(\epsilon^{-2} \log(b\epsilon^{-4} \log^2 m))$

87 / 250

问题3：点查询

流模型 (The Streaming Model)

A stream:

$$\sigma = \langle a_1, a_2, \dots, a_n \rangle, a_i \in [m],$$

then we can define a frequency vector

$$\mathbf{f} = (f_1, \dots, f_m), \text{ s.t. } \sum_{1 \leq i \leq m} f_i = n.$$

定义1.18[点查询 (point query) 问题]

给定数据流 σ 和事先未知查询 a_i , 输出 f_i 。

88 / 250

问题3：点查询

$$\ell_p \text{ 范数: } \|x\|_p = (\sum_i |x_i|^p)^{1/p}$$

定义1.19[ℓ_p -近似点查询 (频度估计)]

给定数据流 σ 以及事先未知查询 a_i 输出 \hat{f}_i 满足 $\hat{f}_i \in f_i \pm \epsilon \|\mathbf{f}\|_p$ 。

- ▷ $\|x\|_1 \geq \|x\|_2 \geq \dots \geq \|x\|_\infty$, p 越大, 估计越准确
- ▷ $\|x\|_0$ 是不同元素的数目
- ▷ $\|x\|_1$ 是流的长度
- ▷ $\|x\|_\infty$ 是最大频度

89 / 250

问题3：点查询- ℓ_1 近似

Misra-Gries 算法[1982]

- ```
/** Maintain a set A consists of pairs (i, f_i) */
1. $A \leftarrow \emptyset$;
2. 对每一个数据流中的元素 e ,
 (a) 如果 $e \in A$, 令 $(e, \hat{f}_e) \leftarrow (e, \hat{f}_e + 1)$
 (b) 否则, 如果 $|A| < 1/\epsilon$, 将 $(e, 1)$ 插入 A
 (c) 否则, 将所有 A 中计数减1
 $((j, \hat{f}_j) \leftarrow (j, \hat{f}_j - 1))$,
 如果 $\hat{f}_j = 0$, 从 A 中删除 $(j, 0)$
3. 对于查询 i , 如果 $i \in A$, return \hat{f}_i , 否则 return 0;
```

91 / 250

## 问题3：点查询- $\ell_1$ 近似

$$\|x\|_1 = n$$

定理1.24

Misra-Gries算法空间代价为  $O(\epsilon^{-1} \log mn)$ , 对任意查询  $i$ , 返回  $\hat{f}_i$  满足

$$f_i - \epsilon n \leq \hat{f}_i \leq f_i$$

- ▷ 如果不发生减1的情况, 那么  $\hat{f}_i = f_i$
- ▷ 如果发生了减1的情况, 有  $\hat{f}_i < f_i$
- ▷ 假设共发生了  $c$  次减1的情况, 总数减少  $\frac{c}{\epsilon} \leq n$
- ▷ 每个计数至多减少  $c$ ,  $\hat{f}_i \geq f_i - c \geq f_i - \epsilon n$

92 / 250

## 问题3：点查询- $\ell_1$ 近似

### Metwally 算法[2005]

- ```
/** Maintain a set A consists of pairs (i, f_i) */
1.  $A \leftarrow \emptyset$ ;
2. 对每一个数据流中的元素  $e$ ,
   (a) 如果  $e \in A$ , 令  $(e, \hat{f}_e) \leftarrow (e, \hat{f}_e + 1)$ 
   (b) 否则, 如果  $|A| < 1/\epsilon$ , 将  $(e, 1)$  插入  $A$ 
       否则, 将  $(e, \text{MIN} + 1)$  插入  $A$ , 并删除
       一个满足  $\hat{f}_e = \text{MIN}$  的元素
       /**这里,  $\text{MIN} = \min\{\hat{f}_e\}$  */
3. 查询  $i$ , 如果  $i \in A$ , return  $\hat{f}_i$ , 否则 return MIN;
```

93 / 250

问题3：点查询- ℓ_1 近似

定理1.25

Metwally算法空间代价为 $O(\epsilon^{-1} \log mn)$, 对任意查询 i , 返回 \hat{f}_i 满足

$$f_i \leq \hat{f}_i \leq f_i + \epsilon n$$

- ▷ 如果不发生删除的情况, 那么 $\hat{f}_i = f_i$
- ▷ 如果删除, 计数一定不大于删除后的MIN, 有 $\hat{f}_i \geq f_i$
- ▷ A 中元素和总是 n , $\text{MIN} \leq n \Rightarrow \text{MIN} \leq \epsilon n$
- ▷ 每个元素至多超出真实值MIN, $\hat{f}_i \leq f_i + \epsilon n$

94 / 250

问题3：点查询- ℓ_1 近似

The (turnstile) Streaming Model

A stream including deletion:

$$\sigma = \langle \pm a_1, \pm a_2, \dots, \pm a_n \rangle, a_i \in [m],$$

then we can define a frequency vector

$$\mathbf{f} = (f_1, \dots, f_m).$$

- ▷ 初始化向量 $\mathbf{f} \in \mathbb{R}^m$ 为0
- ▷ 每次更新 (i, Δ) , 使得 $f_i \leftarrow f_i + \Delta$, f_i 可以比0小
- ▷ $\Delta = 1$ 即是之前所用的流模型 (Vanilla)
- ▷ $\Delta > 0$: 流模型 (Cash Register)

96 / 250

问题3：点查询- ℓ_1 近似

之前讲述方法有局限性

- ▷ 不支持删除操作
令 $\Delta = \pm 1$ 即可定义插入、删除操作
- ▷ Misra-Gries和Metwally算法不支持数据流连接‘o’
无法简单地根据 σ_1 和 σ_2 的结果计算 $\sigma_1 \circ \sigma_2$ 的结果

定义1.20[Sketches]

定义在流数据 σ 上的数据结构 $DS(\sigma)$ 是一个Sketch, 如果存在一个Space-Efficient的合并算法COMB使得

$$\text{COMB}(DS(\sigma_1), DS(\sigma_2)) = DS(\sigma_1 \circ \sigma_2).$$

定义1.21[Linear Sketches]

定义在 $[n]$ 上的流数据 σ 上的sketching算法输出 $\text{sk}(\sigma)$, 如果 $\text{sk}(\sigma)$ 取值为维度 $\ell = \ell(n)$ 的向量, 并且是 $\mathbf{f}(\sigma)$ 的线性函数, 那么 $\text{sk}(\sigma)$ 是一个linear sketch, ℓ 是sketch的维度。

97 / 250

问题3：点查询- ℓ_1 近似

Count-Min Sketch 算法[2005]

1. $k = \frac{4}{\epsilon}$;
2. $C[1 \dots k] \leftarrow 0$;
3. 随机选择一个2-wise独立哈希函数 $h: [m] \rightarrow [k]$;
4. **for** 每一次更新 (j, c) **do**
5. $C[h(j)] = C[h(j)] + c$;
6. **return** $\hat{f}_a = C[h(a)]$ 对查询元素 a ;

- ▷ 空间代价 $O(\frac{1}{\epsilon} \cdot \log n)$
- ▷ 针对Cash Register模型, 有近似质量保证
- ▷ 正确性: $\hat{f}_a \geq f_a$ 并且 $\Pr[\hat{f}_a \geq f_a + \epsilon \|\mathbf{f}_{-a}\|_1] \leq \frac{1}{4}$

98 / 250

问题3：点查询- ℓ_1 近似

Count-Min+ Sketch 算法[2005]

1. $C[1 \dots t][1 \dots k] \leftarrow 0$, $k = \frac{4}{\epsilon}$ and $t = \lceil \log \frac{1}{\delta} \rceil$;
2. 随机选择 t 个2-wise独立哈希函数 $h_i: [m] \rightarrow [k]$;
3. **for** 每一次更新 (j, c) **do**
4. **for** $i = 1$ **to** t **do**
5. $C[i][h_i(j)] = C[i][h_i(j)] + c$;
6. **return** $\hat{f}_a = \min_{1 \leq i \leq t} C[i][h_i(a)]$ 对查询元素 a ;

- ▷ 空间代价 $O(\frac{1}{\epsilon} \log \frac{1}{\delta} \cdot \log n)$
- ▷ 针对Cash Register模型, 有近似质量保证

99 / 250

问题3：点查询- ℓ_1 近似

定理1.26[Count-Min+]

Count-Min+ 算法以 $1 - \delta$ 概率给出 ℓ_1 近似-点查询的 $(1 + \epsilon)$ 近似, 具体地

$$f_a \leq \hat{f}_a \leq f_a + \epsilon \|\mathbf{f}_{-a}\|_1$$

这里, $\|\mathbf{f}_{-a}\|_1 = \|\mathbf{f}\|_1 - f_a$

100 / 250

问题3：点查询- ℓ_1 近似

- ▷ 考虑更一般的Turnstile模型, 即 (j, c) 中 c 可以为负数
- ▷ 使用Count-Min算法, $\Pr[\hat{f}_a \in f_a \pm \epsilon \|\mathbf{f}_{-a}\|_1] \leq \frac{1}{4}$
- ▷ 利用median技术, 可以将概率改进至 $1 - \delta$

Count-Min++ Sketch 算法[2005]

1. $C[1 \dots t][1 \dots k] \leftarrow 0$, $k = \frac{4}{\epsilon}$ and $t = \lceil \log \frac{1}{\delta} \rceil$;
2. 随机选择 t 个2-wise独立哈希函数 $h_i: [m] \rightarrow [k]$;
3. **for** each update (j, c) **do**
4. **for** $i = 1$ **to** t **do**
5. $C[i][h_i(j)] = C[i][h_i(j)] + c$;
6. **return** $\hat{f}_a = \text{median}_{1 \leq i \leq t} C[i][h_i(a)]$ 对查询 a ;

- ▷ 空间代价 $O(\epsilon^{-1} \log \delta^{-1} \cdot \log n)$

104 / 250

问题3：点查询- ℓ_1 近似

定理1.27[Count-Min++]

Count-Min++ 算法以 $1 - \delta$ 概率给出 ℓ_1 近似-点查询的 $(1 \pm \epsilon)$ 近似, 具体

$$|\hat{f}_a - f_a| \leq \epsilon \|\mathbf{f}_{-a}\|_1$$

支持删除操作 (Turnstile模型)

106 / 250

问题3：点查询- ℓ_2 近似

质量更好的估计结果

CountSketch 算法[2004]

1. $C[1 \dots k] \leftarrow 0$, $k = \frac{3}{\epsilon^2}$;
2. 随机选择1个2-wise独立哈希函数 $h: [m] \rightarrow [k]$;
3. 随机选择1个2-wise独立哈希函数 $g: [m] \rightarrow \{-1, 1\}$;
4. **for** each update (j, c) **do**
5. $C[h(j)] = C[h(j)] + c \cdot g(j)$;
6. **return** $\hat{f}_a = g(a) \cdot C[h(a)]$ for query a ;

- ▷ 空间代价 $O(\frac{1}{\epsilon^2} \cdot \log n)$
- ▷ 针对Turnstile模型, 以常数概率, 有 $(1 \pm \epsilon)$ 近似比

109 / 250

问题3：点查询- ℓ_2 近似

利用Median技术改进估计质量

CountSketch+ 算法

1. $C[1 \dots t][1 \dots k] \leftarrow \mathbf{0}, k = \frac{3}{2}, t = O(\log \frac{1}{\delta});$
2. 随机选择 t 个 2-wise 独立哈希函数 $h_i: [m] \rightarrow [k];$
3. 随机选择 t 个 2-wise 独立哈希函数 $g_i: [m] \rightarrow \{-1, 1\};$
4. **for** each update (j, c) **do**
5. **for** each $i = 1$ to t **do**
6. $C[i][h_i(j)] = C[i][h_i(j)] + c \cdot g_i(j);$
7. **return** $\hat{f}_a = \text{median}_{1 \leq i \leq t} g_i(a) C[i][h_i(a)]$ for query $a;$

空间代价 $O(\frac{1}{\epsilon^2} \log \frac{1}{\delta} \cdot \log n)$

113 / 250

问题3：点查询- ℓ_2 近似

$t = \log \frac{1}{\delta}$ 使得切尔诺夫不等式

定理 1.28 [CountSketch+]

CountSketch+ 算法至少以 $1 - \delta$ 概率给出 ℓ_2 近似-点查询的 $(1 \pm \epsilon)$ 近似估计，具体地

$$|\hat{f}_a - f_a| \leq \epsilon \|f_{-a}\|_2$$

114 / 250

问题3：点查询

点查询（频度估计）问题的算法

Method	$\hat{f}_a - f_a \in$	Space	Pr	Mod.
Misra-Gries	$[-\epsilon \ f_{-a}\ _1, 0]$	$O(\frac{\log mn}{\epsilon}) = O(\frac{\log n}{\epsilon})$	0	+
Metwally	$[0, \epsilon \ f_{-a}\ _1]$	$O(\frac{\log mn}{\epsilon}) = O(\frac{\log n}{\epsilon})$	0	+
CountSketch+	$\pm \epsilon \ f_{-a}\ _2$	$O(\frac{1}{\delta^2} \log \frac{1}{\delta} \log n)$	δ	\pm
Count-Min+	$[0, \epsilon \ f_{-a}\ _1]$	$O(\frac{1}{\epsilon} \log \frac{1}{\delta} \log n)$	δ	+
Count-Min++	$\pm \epsilon \ f_{-a}\ _1$	$O(\frac{1}{\epsilon} \log \frac{1}{\delta} \log n)$	δ	\pm

▷ $m < n$

▷ + 表示 Cash Register Model, \pm 表示 Turnstile Model

115 / 250

问题4：频度矩估计

假设我们现在处理 Vanilla Model

定义 1.22 [频度矩--Frequency Moments]

给定数据流 $\sigma = \langle a_1, \dots, a_n \rangle$, 假设 $a_i \in [m]$, 令 $\mathbf{f} = \mathbf{f}(\sigma) = (f_1, f_2, \dots, f_m)$, 显然, $\sum f_i = n$. 数据流 σ 的第 k 阶频度矩 (kth frequency moment) 表示为 $F_k(\sigma)$ 或者 F_k ,

$$F_k = \sum_{i=1}^m f_i^k = \|\mathbf{f}\|_k^k$$

事实上, 这里 k 可以扩展为实数 $k \geq 0$, 很多方法可以扩展到 Cash Register Model

116 / 250

问题4：频度矩估计

- ▷ F_0 是不重复元素的数目
- ▷ F_1 是计数问题
- ▷ F_2 可表示关系数据自连接 (self join) 操作结果大小
- ▷ 这部分给出 AMS 算法, 为 $F_k (k \geq 2)$ 估计问题提供亚线性空间的算法

117 / 250

问题4：频度矩估计

Basic AMS 算法

1. $(L, r, a) \leftarrow (0, 0, 0);$
2. **for** 第 j 个元素 a_j **do**
3. $L \leftarrow L + 1;$
4. $\beta \leftarrow$ random bit with $\Pr[\beta = 1] = \frac{1}{L};$
5. **if** $\beta = 1$ **then**
6. $a \leftarrow a_j;$
7. $r \leftarrow 0;$
8. **if** $a_j == a$ **then**
9. $r \leftarrow r + 1;$
10. **return** $X = L(r^k - (r - 1)^k);$

118 / 250

问题4：频度矩估计

空间代价

- ▷ 存储 L 和 r 需要 $\log n$ 位
- ▷ 存储 a 需要 $\log m$ 位
- ▷ 共计 $O(\log m + \log n)$ 位代价

直观思想

- ▷ 随机选取一个位置 $y \in [n]$
- ▷ $a = a_y$
- ▷ $r = |\{j | j \geq y \text{ 并且 } a_j = a_y\}|$

$$\Pr[|X - \mathbf{E}[X]| \geq \epsilon \mathbf{E}[X]] \leq \frac{kn^{1-\frac{1}{k}} F_k^2}{\epsilon^2 F_k^2} = \frac{kn^{1-\frac{1}{k}}}{\epsilon^2}$$

119 / 250

问题4：频度矩估计

- ▷ Basic AMS 算法的方差太大, 无法直接利用 Median 技术将获得 $(1 + \epsilon)$ 近似的失败概率降低到满足要求
 - ▷ 首先需要利用 averaging 技术, 将方差降低
- 更一般化的表达

定理 1.29 [Median-of-Means]

存在常数 c 满足: 令 X 是一个实数量 Q 的无偏估计, 令 $\{X_{ij}\}_{i \in [t], j \in [k]}$ 是与 X 同分布的独立变量集合,

$$t = c \log \frac{1}{\delta} \text{ 并且 } b = \frac{3 \text{var}[X]}{\epsilon^2 \mathbf{E}[X]^2}$$

令 $Z = \text{median}_{i \in [t]} (\frac{1}{b} \sum_{j=1}^b X_{ij})$, 有 $\Pr[|Z - Q| \geq \epsilon Q] \leq \delta$.

124 / 250

问题4：频度矩估计

Final AMS 算法

1. 利用Median-of-Mean技术调用Basic AMS算法;
2. 计算 $t = c \log \frac{1}{\delta}$ 个平均值;
3. 每个平均值是 $r = \frac{3k}{\epsilon^2} \cdot n^{1-\frac{1}{k}}$ 次调用的平均值;
4. **return** t 个数值的中间值;

$$\triangleright r = \frac{3\text{var}[X]}{\epsilon^2 \mathbb{E}[X]^2} \leq \frac{3kn^{1-\frac{1}{k}}F_k^2}{\epsilon^2 F_k^2} = \frac{3k}{\epsilon^2} \cdot n^{1-\frac{1}{k}}$$

- ▷ 空间代价：
 $tr \cdot O(\log m + \log n) = O(\frac{1}{\epsilon^2} \log \frac{1}{\delta} \cdot kn^{1-\frac{1}{k}} (\log m + \log n))$
- ▷ 忽略 $\log m, \log n, \log \frac{1}{\delta}$, 对常数 k , 有 $\tilde{O}(\frac{1}{\epsilon^2} n^{1-\frac{1}{k}})$
- ▷ 更优的: $\tilde{O}(\frac{1}{\epsilon^2} n^{1-\frac{2}{k}})$, 下界 $\tilde{\Omega}(n^{1-\frac{2}{k}})$ 和 $\tilde{\Omega}(\epsilon^{-2})$

126 / 250

问题4：频度矩估计

上述算法在 $k=2$ 时, 代价为 $\tilde{O}(\epsilon^{-2} n^{0.5})$, 过大

Basic F_2 AMS 算法(Tug-of-War Sketch)

1. 随机选择4-wise独立的哈希函数 $h: [m] \rightarrow \{-1, 1\}$;
2. $x \leftarrow 0$;
3. **for** each update (j, c) **do**
4. $x \leftarrow x + c \cdot h(j)$;
5. **return** x^2 ;

▷ 空间代价: $O(\log n)$

▷ $\mathbb{E}[X^2] = F_2$

▷ $\text{var}[X^2] \leq 2F_2^2$

127 / 250

问题4：频度矩估计

- ▷ 利用media-of-mean技术, 可获得 (ϵ, δ) 近似算法
- ▷ 令 $t = c \log \frac{1}{\delta}$
- ▷ 令 $k = \frac{3\text{var}[X^2]}{\epsilon^2 \mathbb{E}[X^2]^2} \leq \frac{6}{\epsilon^2}$
- ▷ 空间代价 $\tilde{O}(\epsilon^{-2} \log \frac{1}{\delta}) = O(\frac{1}{\epsilon^2} \log \frac{1}{\delta} \log n)$

130 / 250

问题5：频繁元素

定义1.23[频繁元素, Frequent Items, heavy hitters]

给定数据流 σ 和参数 k , 输出频繁元素集合 $\{j \mid f_j > n/k\}$.

利用点查询的算法求解

- ▷ 令 $\epsilon = 1/k$
 - ▷ 执行Misra-Gries 算法一次, 获得 A
 - ▷ 输出 A 中元素
- 如果 $f_j > n/k$, 那么 $\tilde{f}_j > 0$, $j \in A$

131 / 250

问题5：频繁元素

- ▷ 想估计最大频度, 是否可以在亚线性代价内实现?
- ▷ 考虑一个hard情况, 每个元素频度都差不多, 最后一个元素决定最大频度是多少
- ▷ 放松要求: 最频繁的元素是heavy的, 设计算法

定义1.24[Heavy Hitters]

给定参数 $\phi \in (0, 1)$ 和 $p \in [1, \infty)$, 对于数据流 σ , 假设对应的频度向量为 \mathbf{f} , heavy hitters定义如下

$$HH_\phi^p(\mathbf{f}) = \{i \in [m] : f_i \geq \phi \|\mathbf{f}\|_p\}$$

我们关注 $p=1$ 的情况

定义1.25[ϕ -heavy]

元素 i 是 ϕ -heavy的, 如果 $f_i \geq \phi \sum f_j$

132 / 250

问题5：频繁元素

定义1.26[Approximate Heavy Hitters for $p=1$]

给定参数 $\phi, \epsilon \in (0, 1)$, 对于数据流 σ , 假设对应的频度向量为 \mathbf{f} , 计算一个集合 $S \subseteq [m]$ 使得

$$HH_\phi \subseteq S \subseteq HH_{(1-\epsilon)\phi}$$

利用 ℓ_1 点查询算法求解, 设计 (ϵ', δ') 近似算法

- ▷ 利用Count-Min++算法实现 (ϵ, δ) 近似点查询算法
 - ▷ 令 $\epsilon = \epsilon' \phi / 2$, $\delta = \delta' / m$, 为每一个元素 $i \in [m]$ 估计 \hat{f}_i
 - ▷ 返回集合 $S = \{i \in [m] : \hat{f}_i \geq (\phi - \epsilon' \phi / 2) \|\mathbf{f}\|_1\}$
- 插入更新 $\|\mathbf{f}\|_1$ 非常容易获得, 有删除就很难
- ▷ 算法的空间代价: $O(\frac{1}{\phi \epsilon'} (\log \frac{1}{\delta'} + \log m) \log n)$

133 / 250

问题6：固定大小采样

- ▷ 问题描述: 时刻维护一个大小为 s 的采样, 使得 σ 中的每个元素出现的概率均是 $\frac{s}{n}$

水库抽样算法

1. $n \leftarrow 0$;
2. 用流数据的前 s 个元素初始化 $A[1, \dots, s]$, $n \leftarrow s$;
3. **for** 每一个元素 x **do**
4. x 以 $\frac{s}{n+1}$ 概率进入 A ;
5. **if** x 加入 A **then**
6. 任意选择 A 中元素剔除;
7. $n \leftarrow n + 1$;

136 / 250

问题7：Bloom Filter

定义1.27[Membership Problem]

令 U 是整数集合 $\{1, \dots, |U|\}$, S 是 U 的一个大小为 n 的子集, 预处理 S , 给定整数 $q \in U$, 判定是否 $q \in S$.

- ▷ U 中整数可以用 $w = \log |U|$ 位表示
- ▷ 哈希表, 查询期望时间是 $O(1)$, 总代价 $O(nw)$ 位
- ▷ 是否可以在 $o(nw)$ 代价解决?
如果要求精确答案, 那么至少用 $\log \binom{|U|}{n}$ 位, 当 $|U| \gg n$, 即 $\Omega(nw)$
- ▷ 放松要求
允许false positives, 不允许false negatives

138 / 250

问题7：Bloom Filter

定义1.28[Approximate Membership Problem]

给定整数 $q \in U$,

- ▷ $q \in S \Rightarrow$ 输出 ‘yes’
- ▷ $q \notin S \Rightarrow$ 以至少 $1 - \delta$ 概率输出 ‘no’

139 / 250

问题7：Bloom Filter

近似哈希的方法

1. 令 \mathcal{H} 是一族通用哈希函数: $|U| \rightarrow [m]$, $m = \frac{n}{\delta}$;
2. 随机选取 $h \in \mathcal{H}$, 并维护数组 $A[m]$;
3. **for** $i \in S$ **do**
4. $A[h(i)] = 1$;
5. 给定查询 q , **return** ‘yes’ 当且仅当 $A[h(i)] = 1$;

- ▷ $q \in S \Rightarrow$ 输出 ‘yes’
- ▷ $q \notin S \Rightarrow$

$$\sum_{j \in S} \Pr[h(q) = h(j)] \leq \frac{n}{m} = \delta$$

140 / 250

问题7：Bloom Filter

定义1.29[Ideal Hash Function]

哈希函数 $h: U \rightarrow [m]$ 是理想的 (ideal), 如果对每一个 $k \in U$, $h(k)$ 均匀独立地从 $[1, m]$ 间取值。

Bloom Filter方法

1. 令 \mathcal{H} 是一族独立的理想哈希函数: $|U| \rightarrow [m]$;
2. 随机选取 $h_1 \dots h_d \in \mathcal{H}$, 并维护数组 $A[m]$;
3. **for** $i \in S$ **do**
4. **for** $j \in [1, d]$ **do**
5. $A[h_j(i)] = 1$;
6. 给定查询 q ,
 return ‘yes’ 当且仅当 $\forall j \in [d], A[h_j(q)] = 1$;

代价: $m = O(n \log \frac{1}{\delta}) \Rightarrow O(n \log \frac{1}{\delta})$

141 / 250

▣ 概率基础

▣ 亚线性空间算法

▣ 亚线性时间算法

▣ 并行模型算法

145 / 250

亚线性时间算法

高效算法的标准是什么?

- ▷ 多项式时间算法
- ▷ 线性时间算法
 - 输入大小为 n , 算法的时间代价是 cn

如果时间资源受限, 甚至不允许完整读取数据一遍, 我们能做什么?

- ▷ 无法回答 “for all” 或者 “exactly” 类似的问题: 数组中是否所有元素均比10大? 数组中比10大的有多少?
- ▷ 或许能够回答下列问题: 数组中比10大的元素大概有多少? 数组中的平均数是否以很大概率大于10?

算法的计算目标需要修改

146 / 250

亚线性时间算法

算法的计算目标需要修改

- ▷ 对于绝大多数问题, 我们只能给出近似的答案 什么是 “近似”
- ▷ 经典的近似算法概念:
 - 算法的结果有一个度量函数: 路径长度
 - 算法结果度量值与最优解度量值 “接近”: 近似比
- ▷ Property Testing (判定问题)
 - 算法的结果是 yes 或者 no
 - 算法输出要么与输入 I 对应的答案相同, 要么与 I' 对应的答案相同, 这里 I' 与 I 非常接近

147 / 250

问题1：点集合的直径

点集合的直径问题

输入: m 个点, 距离用矩阵 D 表示

- (1) D_{ij} 是 i 点到 j 点的距离
- (2) D 是对称的, 并且满足三角不等式
 $D_{ij} \leq D_{ik} + D_{kj}$

输出: (i, j) 使得 D_{ij} 是最大的, D_{ij} 是点集合的直径

148 / 250

问题1：点集合的直径

亚线性求解直径算法 (Indyk's Algorithm)

*/** input D */*

1. 选取任意一个 k ;
2. 选取 l 最大化 D_{kl} ;
3. **return** (k, l)

近似比分析: $opt/2 \leq D_{kl} \leq opt$

$$\begin{aligned} opt = D_{ij} &\leq D_{ik} + D_{kj} && \text{(三角不等式)} \\ &\leq D_{kl} + D_{kl} && (l \text{ 是最大化选择}) \\ &\leq 2D_{kl} \end{aligned}$$

运行时间: $O(m) = O(\sqrt{n})$

149 / 250

问题2：连通分量的数目

连通分量的数目 (#CC)

输入: $G = (V, E)$, ϵ , $d = \deg(G)$

图G用邻接链表表示

$|V| = n$, $|E| = m \leq dn$

输出: y , 令 C 为 #CC

$C - \epsilon n \leq y \leq C + \epsilon n$ (additive apprx.)

#CC可以在线性时间求解 (DFS或者BFS)

150 / 250

问题2：连通分量的数目

n_v : 顶点 v 所属的连接分量中的节点数目

$A: A \subseteq V$ 是一个连通分量的点集合

$$\sum_{u \in A} \frac{1}{n_u} = \sum_{u \in A} \frac{1}{|A|} = 1$$

$$C = \#CC = \sum_{u \in V} \frac{1}{n_u}$$

为什么用这种表示?

- ▷ 计算 $\frac{1}{n_u}$ 需要 $O(n)$ 时间?
- ▷ 需要对 $O(n)$ 项求和?
- ▷ 可以 estimate

151 / 250

问题2：连通分量的数目

估计 $C = \sum_{u \in V} \frac{1}{n_u} \Rightarrow$ 估计 $\frac{1}{n_u} \Rightarrow$ 估计 $\sum_{u \in V} \frac{1}{n_u}$

估计 $\frac{1}{n_u}$

想法: n_u 很大, 精确计算很难, 但此时 $\frac{1}{n_u}$ 很小, 可以用一个很小的常量代替 $\frac{1}{n_u}$ (0 或者 $\epsilon/2$)

$$\hat{n}_u = \min\{n_u, \frac{2}{\epsilon}\}$$

$$\hat{C} = \sum_{u \in V} \frac{1}{\hat{n}_u}$$

引理 1.6

$\forall u \in V$, 有 $|\frac{1}{n_u} - \frac{1}{\hat{n}_u}| \leq \epsilon/2$, 即 $|C - \hat{C}| \leq \frac{\epsilon n}{2}$

152 / 250

问题2：连通分量的数目

估计 $C = \sum_{u \in V} \frac{1}{n_u} \Rightarrow$ 估计 $\frac{1}{n_u} \Rightarrow$ 估计 $\sum_{u \in V} \frac{1}{n_u}$

$$\hat{n}_u = \min\{n_u, \frac{2}{\epsilon}\}; \hat{C} = \sum_{u \in V} \frac{1}{\hat{n}_u}$$

计算 \hat{n}_u 算法

1. 从 u 开始做先广遍历 (BFS);
2. **while** BFS 访问过的节点数量 $\leq \frac{2}{\epsilon}$ **do**
3. 继续做 BFS 遍历;
4. **if** 没有新的节点可以遍历 **then**
5. **return** BFS 访问过的节点数量;
6. **return** $2/\epsilon$;

运行时间: $O(d \cdot 1/\epsilon)$

153 / 250

问题2：连通分量的数目

估计 $C = \sum_{u \in V} \frac{1}{n_u} \Rightarrow$ 估计 $\frac{1}{n_u} \Rightarrow$ 估计 $\sum_{u \in V} \frac{1}{n_u}$

$$\hat{n}_u = \min\{n_u, \frac{2}{\epsilon}\}; \hat{C} = \sum_{u \in V} \frac{1}{\hat{n}_u}$$

亚线性连通分量数目求解算法 (用 \hat{C} 估计 \hat{C})

1. $r \leftarrow b/\epsilon^2$;
2. 随机从 V 中选取 $U = \{u_1, \dots, u_r\}$;
3. 计算所有的 \hat{n}_{u_i} ;
4. **return** $\tilde{C} = \frac{n}{r} \sum_{u \in U} \frac{1}{\hat{n}_{u_i}}$;

运行时间: $O(d \cdot 1/\epsilon \cdot 1/\epsilon^2) = O(d/\epsilon^3) = o(|G|)$

近似性能: $r = O(\frac{1}{\epsilon^2} \log \frac{1}{\delta}) \Rightarrow \Pr[|C - \hat{C}| \geq \epsilon n] \leq \delta$

154 / 250

问题3：近似最小支撑树

最小支撑树 (Min Spanning Tree)

输入: $G = (V, E)$, ϵ , $d = \deg(G)$

图G用邻接链表表示

边 (u, v) 的权重是 $w_{uv} \in \{1, 2, \dots, w\} \cup \{\infty\}$

输出: \hat{M} , 令 M 为 $\min_{T \text{ spans } G} \{W(T)\}$

$$(1 - \epsilon)M \leq \hat{M} \leq (1 + \epsilon)M$$

MST问题可以在多项式时间求解, 例如Kruskal算法

158 / 250

问题3：近似最小支撑树

我们需要重新利用分解的方式定义 M

$G^{(i)} = (V, E^{(i)})$, 这里 $E^{(i)} = \{(u, v) | w_{uv} \leq i\}$

$C^{(i)} = \#CC \text{ in } G^{(i)}$

让我考虑两个例子

- ▷ $w = 1$: 只有权重为1的边, 且连通
 $M = n - 1$
- ▷ $w = 2$: 有权重为1和2的边, 且连通
 $M = n - 1 + C^{(1)} - 1 = n - 2 + C^{(1)}$

159 / 250

问题3：近似最小支撑树

定理 1.31 [MST]

$$M = n - w + \sum_{i=1}^{w-1} C^{(i)}$$

α_i : 任一MST中权重为 i 的边的数目

$$\sum_{i>1} \alpha_i = C^{(1)} - 1$$

$$M = \sum_{i=1}^w i \cdot \alpha_i = \sum_{i=1}^w \alpha_i + \sum_{i=2}^w \alpha_i + \dots + \sum_{i=w}^w \alpha_i$$

$$= C^{(0)} - 1 + C^{(1)} - 1 + \dots + C^{(w-1)} - 1$$

$$= n - 1 + C^{(1)} - 1 + \dots + C^{(w-1)} - 1 = n - w + \sum_{i=1}^{w-1} C^{(i)}$$

160 / 250

问题3：近似最小支撑树

利用连通分量数目估计的算法来估计MST大小

$$\tilde{O}(d/\epsilon^3) \Rightarrow \Pr[|C - \hat{C}| \leq \epsilon n] \geq 1 - \delta$$

亚线性近似最小支撑树算法（计算 \hat{M} ）

1. **for** $i = 1$ to $w - 1$ **do**
2. $\hat{C}^{(i)} = \text{appro. \#CC of } G^{(i)} \text{ within } (\epsilon' = \frac{\epsilon}{2w}) \cdot n;$
with probability $\geq 1 - \delta' = 1 - \frac{\delta}{w}$
3. **return** $\hat{M} = n - w + \sum_{i=1}^{w-1} \hat{C}^{(i)};$

单次时间: $\tilde{O}(d \cdot 1/\epsilon'^3) = \tilde{O}(dw^3/\epsilon^3)$

共计时间: $\tilde{O}(dw^4/\epsilon^3) = o(|G|)$

近似性能: $\Pr[|\hat{M} - M| \geq \epsilon M] \leq \delta$

161 / 250

问题4：计算图的平均度

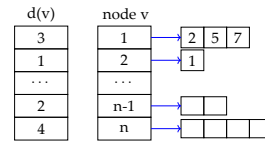
图的平均度（Average Degree）

输入: $G = (V, E)$

图G用邻接链表表示

输出: 平均度 $\bar{d} = \frac{\sum_{u \in V} d(u)}{n}$

假设: (1) G是简单图（无平行边、自环）， $\Omega(n)$ 条边
(2) G由邻接链表和存储度数的数组表示



▷ degree query: $d(v)$

▷ neighbor query: $N(v, j)$
v的第j个邻居节点

164 / 250

问题4：计算图的平均度

利用Naive Sampling的方法

- ▷ 随机选取s个节点 $\{v_1, \dots, v_s\}$, 返回 $\frac{1}{s} \sum_{i=1}^s d(v_i)$
- 简单的下界分析
- ▷ 空图的平均度数为 $\bar{d} = 0$, 加一条边后变为 $\bar{d} = 2/n$.
 - 为了区分上述两种情况, 基本上需要访问所有的节点, 即 $\Omega(n)$.
- ▷ 令 G_1 为大小为 n 的环: $\bar{d} = 2$;
令 G_2 由 $(n - c\sqrt{n})$ -环和 $c\sqrt{n}$ -团组成: $\bar{d} \approx 2 + c^2$
 - 大致需要 $\frac{n}{c\sqrt{n}} = \frac{\sqrt{n}}{c}$ 步找到这个团, 即 $\Omega(\sqrt{n})$

165 / 250

问题4：计算图的平均度

算法思想

- ▷ 将具有相似或者相同度数的节点分组
- ▷ 估计每个分组的平均度数, 分组的方差都有界
 - 对于任意数值数组这个方法并不适用

分桶策略: $\beta = \frac{\epsilon}{c_\beta};$

$$t = \lceil \log_{(1+\beta)} n \rceil + 1 = O\left(\frac{\log n}{\log(1+\beta)}\right) \approx O\left(\frac{\log n}{\beta}\right) = O\left(\frac{\log n}{\epsilon}\right)$$

- ▷ $B_i = \{v | (1+\beta)^{i-1} < d(v) \leq (1+\beta)^i\}, 0 \leq i \leq t-1$
- ▷ B_i 的总度数: $(1+\beta)^{i-1}|B_i| < d(B_i) \leq (1+\beta)^i|B_i|$
- ▷ G 的总度数: $\sum_{i=0}^{t-1} (1+\beta)^{i-1}|B_i| < \sum d(u) \leq \sum_{i=0}^{t-1} (1+\beta)^i|B_i|$
- ▷ $\frac{\sum_{i=0}^{t-1} (1+\beta)^{i-1}|B_i|}{n} < \bar{d} \leq \frac{\sum_{i=0}^{t-1} (1+\beta)^i|B_i|}{n}$ 重点是估计 $|B_i|/n$

166 / 250

问题4：计算图的平均度

第一个算法（naive sampling）

平均度估计算法I

1. 从 V 中选取 S 个样本点;
2. $S_i \leftarrow S \cap B_i;$ //利用 S_i 来估计 B_i
3. $\rho_i \leftarrow \frac{|S_i|}{|S|};$
4. **return** $\hat{\bar{d}} = \sum_{i=0}^{t-1} \rho_i (1+\beta)^{i-1};$ //用下界, 低估

- ▷ 令 $X_j^i = 1$ 表示样本 s_j 落到桶 B_i 中; 否则, $X_j^i = 0$

$$\triangleright \mathbf{E}[\rho_i] = \mathbf{E}\left[\frac{|S_i|}{|S|}\right] = \frac{\mathbf{E}[\sum_{j=1}^{|S|} X_j^i]}{|S|} = \frac{|S_i| \cdot |B_i|}{|S| \cdot n} = \frac{|B_i|}{n}$$

167 / 250

问题4：计算图的平均度

算法I的问题: 较小的 B_i 的估计效果会很差

考虑完全二分图 $K_{3,n-3}$

- ▷ 有3个节点度为 $n-3$, 有 $n-3$ 个节点度为3
- ▷ $B_a: (1+\beta)^{a-1} < 3 \leq (1+\beta)^a$ $|B_a| = n-3$
- ▷ $B_b: (1+\beta)^{b-1} < n-3 \leq (1+\beta)^b$ $|B_b| = 3$
- ▷ B_a 和 B_b 在 G 中连接的边数相同, 均是 $3(n-3)$
- ▷ $\bar{d} = \frac{3(n-3) + (n-3)3}{n} = 6 - \frac{18}{n} \approx 6$
- ▷ 算法I很可能未采样到 B_b 否则需 $\Omega(n)$ 期望时间
- ▷ $\hat{\bar{d}} = 3 \Rightarrow \frac{\bar{d}}{\hat{\bar{d}}}$ 至少为2（并没有分析其它worst cases）

168 / 250

问题4：计算图的平均度

改进算法（希望获得稳定的近似比为2的算法）

对于较小的桶置0, 假定已知一个 \bar{d} 的下界 α

平均度估计算法II

1. 从 V 中采样 $S, |S| = \Theta(\sqrt{n/\alpha} \cdot \text{poly}(\log n, 1/\epsilon));$
2. $S_i \leftarrow S \cap B_i;$
3. **for** $i \in \{0, \dots, t-1\}$ **do**
4. **if** $|S_i| \geq \theta_\rho |S| = \frac{1}{t} \sqrt{\frac{3}{8}} \cdot \frac{\epsilon \alpha}{n} |S|$ **then**
5. $\rho_i \leftarrow \frac{|S_i|}{|S|};$
6. **else**
7. $\rho_i \leftarrow 0;$
8. **return** $\hat{\bar{d}} = \sum_{i=0}^{t-1} \rho_i (1+\beta)^i;$

169 / 250

问题4：计算图的平均度

$$\beta = \frac{\epsilon}{c_\beta}, |S| = \frac{c_s \log n}{\epsilon^{3.5}} \sqrt{\frac{n}{\alpha}} = \tilde{O}\left(\frac{\log n}{\epsilon^{3.5}}\right) \sqrt{\frac{n}{\alpha}},$$

$$t = \lceil \log_{(1+\beta)} n \rceil + 1, |S_i| \geq \frac{1}{t} \sqrt{\frac{3}{8}} \cdot \frac{\epsilon \alpha}{n} |S| = \theta_\rho |S|$$

$$\hat{\bar{d}} = \sum_{i=0}^{t-1} \rho_i (1+\beta)^i \approx \sum_{i=0}^{t-1} \frac{|B_i|}{n} (1+\beta)^i \approx \bar{d}$$

- ▷ 如果对于所有 B_i , ρ_i 都是非常接近 $\frac{|B_i|}{n}$, OK
- ▷ 如果 $|B_i|$ 很大, 根据Chernoff界, 应该OK
- ▷ 如果 $|B_i|$ 很小, 估计的方差可能会很大, 估计为0
 - 本质上是计算边数的2倍, 即 $\bar{d} \cdot n$
 - 当边的一个顶点在小 B_i 中, 误差为2
 - 当边的两个顶点在小 B_i 中, 误差很大, 但这样的边数大小可控

170 / 250

问题4：计算图的平均度

$$\beta = \frac{\epsilon}{c_\beta}, |S| = \frac{c_s \log n}{\epsilon^{3.5}} \sqrt{\frac{n}{\alpha}} = \tilde{O}\left(\frac{\log n}{\epsilon^{3.5}}\right) \sqrt{\frac{n}{\alpha}},$$

$$t = \lceil \log_{(1+\beta)} n \rceil + 1, |S_i| \geq \frac{1}{t} \sqrt{\frac{3}{8} \cdot \frac{c_\alpha}{n}} |S| = \theta_\rho |S|$$

$$\frac{c_s}{c_\beta} = O(\log \frac{t}{\delta}) = O(\log \frac{\log n}{\epsilon \delta}) \Rightarrow \text{至少以 } 1 - \delta_1 - \delta_2 \text{ 概率}$$

$$(0.5 - \epsilon) \bar{d} < \hat{d} < (1 + \epsilon) \bar{d}$$

因此，算法II是一个 $(2 + \epsilon)$ 近似比的算法，时间复杂度为 $\tilde{O}(\sqrt{n} \cdot \text{poly}(\frac{1}{\epsilon}))$

171 / 250

问题4：计算图的平均度

算法II的改进：

- ▷ 将节点分为 U (度数小) 和 $V \setminus U$ (度数大)，近似比中的 2 因为 $E(V \setminus U, U)$ 只计算了一次
- ▷ 估计 $E(V \setminus U, U)$ 的比例，能校正这部分错误
- ▷ 通过从采样得到的节点邻居中随机选取一条边实现
 - 随机邻居查询 (random neighbor query)
 - $rnq(v)$
 - ✓ 获得 $d(v)$
 - ✓ 随机生成 $i \in [1, d(v)]$
 - ✓ 获得 $N(v, i)$

188 / 250

问题4：计算图的平均度

算法II的改进：

$$|B_i| \cdot (1 + \beta)^{i-1} < |E_i| \leq |B_i| \cdot (1 + \beta)^i$$

访问 B_i 中的点来估计 E_i 中 $E(V \setminus U, U)$ 的比例

- ▷ 重复如下操作 m 次
 - 随机选择 $u \in B_i$ 利用 S_i 实现
 - 令 $u' = rnq(u)$ 利用 rnq 实现
 - 如果 $(u, u') \in E(V \setminus U, V \setminus U)$ ，令 $\delta_j = 0$
 - 否则 $(u, u') \in E(V \setminus U, U)$ ，令 $\delta_j = 1$ 如何判断？
- ▷ 返回 $\Delta_i = \frac{\sum \delta_j}{m}$

189 / 250

问题4：计算图的平均度

访问 B_i 中的点来估计 E_i 中 $E(V \setminus U, U)$ 的比例

- ▷ 令 T_i 为 B_i 中 big-small 的边数
- ▷ 简单情况： B_i 中点度数 d 相同

$$\mathbf{E}[\Delta_i] = \Pr[\delta_j = 1] = \frac{T_i}{d \cdot |B_i|}$$

- ▷ 一般情况： B_i 中度数比值在 $(1 + \beta)$ 内

$$\frac{T_i}{|B_i|(1 + \beta)^i} \leq \mathbf{E}[\Delta_i] \leq \frac{T_i}{|B_i|(1 + \beta)^{i-1}}$$

$$|B_i|(1 + \beta)^{i-1} \mathbf{E}[\Delta_i] \leq T_i \leq |B_i|(1 + \beta)^i \mathbf{E}[\Delta_i]$$

- ▷ $\mathbf{E}[\Delta_i]$ 的 $(1 + \epsilon)$ 估计
 $\Rightarrow \Delta_i \rho_i (1 + \beta)^i$ 是 $\frac{T_i}{n}$ 的 $(1 + \epsilon)(1 + \beta)$ 估计

190 / 250

问题4：计算图的平均度

改进算法：近似比 $(1 \pm \epsilon)$ ；运行时间 $\tilde{O}(\sqrt{n} \cdot \text{poly}(1/\epsilon))$

平均度估计算法III

1. Sample S of V , $|S| = \tilde{O}(\frac{L}{\epsilon^2})$, $L = \text{poly}(\frac{\log n}{\epsilon})$, $\theta_\rho = \frac{1}{t} \sqrt{\frac{3\epsilon}{8} \cdot \frac{\alpha}{n}}$;
2. $S_i \leftarrow S \cap B_i$;
3. **for** $i \in \{0, \dots, t-1\}$ **do**
4. **if** $|S_i| \geq \theta_\rho |S|$ **then**
5. $\rho_i \leftarrow \frac{|S_i|}{|S|}$;
6. **for** $v \in S_i$ **do**
7. $u \leftarrow rnq(v)$; $\delta(v) \leftarrow 0$;
8. **if** $u \in B_j \wedge |S_j| < \theta_\rho |S|$ **then** $\delta(v) \leftarrow 1$;
9. $\Delta_i = \sum_{v \in S_i} \delta(v) / |S_i|$;
10. **else**
11. $\rho_i \leftarrow 0$;
12. **return** $\hat{d} = \sum_{i=0}^{t-1} (1 + \Delta_i) \rho_i (1 + \beta)^i$;

191 / 250

问题4：计算图的平均度

消去参数 $\alpha \leq \bar{d}$ ，即度数的下界

平均度估计算法VI

1. $\alpha \leftarrow n$;
2. $\hat{d} \leftarrow -\infty$;
3. **while** $\hat{d} < \alpha$ **do**
4. $\alpha \leftarrow \alpha/2$;
5. **if** $\alpha < \frac{1}{n}$ **then**
6. **return** 0;
7. 利用估计算法 III 在参数 α 下运行，估计结果为 \hat{d} ;
8. **return** \hat{d} ;

近似比： $(1 + \epsilon)$;
 运行时间： $\tilde{O}(\sqrt{n} \cdot \text{poly}(\epsilon^{-1}))$

200 / 250

问题4：计算图的平均度

如果有parallel边或者边带权重，那么估计算法的下界 $\Omega(n)$

- ▷ G_1 : n -cycle graph
- ▷ G_2 : $n - 2$ -cycle graph + one $c \cdot n$ parallel edge
- ▷ $\bar{d}(G_1) = 2$
- ▷ $\bar{d}(G_2) = c + 2$

201 / 250

问题5：Vertex Cover

顶点覆盖 (Vertex Cover)

输入： $G = (V, E)$ 、最大度数 d
 图 G 用邻接链表表示
 $C \subseteq V$ 是一个 VC (Vertex Cover) 当且仅当
 $\forall (u, v) \in E$ 有 $u \in C$ 或者 $v \in C$
 输出：最小 VC 的大小 $|VC|$

- ▷ $|VC| \geq \frac{|E|}{d}$
- ▷ NP-完全问题
- ▷ 存在集中式的 2 近似算法

203 / 250

问题5：Vertex Cover

在亚线性时间内，能有多好的近似算法？

- ▷ 乘性的近似比 ($\frac{|\hat{VC}|}{|VC|}$): No
Graph with 0 edge $\Rightarrow |VC| = 0$
Graph with 1 edge $\Rightarrow |VC| = 1$
区分上述两种情况需要 $\Omega(n)$ 时间
- ▷ 加性的近似比 ($|\hat{VC}| - |VC|$): Hard
乘性近似比的下界是 1.36，很可能是 2

定义 1.30 (α, ϵ) -近似

对于最优解是 y 的最小化问题，如果

$$y \leq \hat{y} \leq \alpha y + \epsilon$$

则称 \hat{y} 是该问题的 (α, ϵ) -近似。

204 / 250

问题5：Vertex Cover

设计想法：利用分布式算法设计亚线性算法

分布式网络

- ▷ max degree d
- ▷ 每个节点是一个处理器，知道它的邻居是谁
- ▷ 同步计算、通信

每一轮（同步）

- ▷ 每个节点计算基于它的输入、随机生成位、通信中收到的信息
- ▷ 向每个邻居节点发送消息
- ▷ 从每个邻居节点收取消息

205 / 250

问题5：Vertex Cover

设计想法：利用分布式算法设计亚线性算法

分布式顶点覆盖（Vertex Cover）问题

- ▷ 输入：网络图即是要计算顶点覆盖的图 G
- ▷ 输出：每个节点知道自己是否属于 VC

Fast Distributed Alg. \Rightarrow SubLinear Time Alg.

- ▷ k 轮分布式算法中，节点 v 最多依赖于距离为 k 的节点，至多 d^k 个
- ▷ 集中式算法可以在 d^k 时间内，模拟分布式算法，计算 v 是否属于 VC

注：如果是随机算法，需要知道所有 d^k 个节点的随机位

206 / 250

问题5：Vertex Cover

- ▷ 存在 fast VC distribute alg (local distributed alg)?
- ▷ 如何利用分布式算法设计亚线性算法

Sublinear VC Alg (Parnas-Ron framework)

1. 独立均一地从 G 中选取 $s = \frac{8}{\epsilon^2}$ 个节点构成 S ;
2. 为每个 $v \in S$ ，构造 k 步邻居导出的子图 $G_k(v)$;
3. 为每个 $v \in S$ ，在 $G_k(v)$ 上模拟分布式算法 \mathcal{D} ;
4. 如果 \mathcal{D} 返回 v 为覆盖节点之一， $X_v = 1$;
5. **return** $\hat{VC} = \frac{n}{s} \cdot \sum_{v \in S} X_v + \frac{\epsilon}{2} n$;

运行时间分析： $O(s \cdot d^k) = O(\frac{d^k}{\epsilon^2})$

207 / 250

问题5：Vertex Cover

- ▷ 存在 fast VC distribute alg (local distributed alg)?

Distributed VC Alg

1. $\tilde{VC} \leftarrow \emptyset$;
2. **for** $i = 1$ to $\log d$ **do**
3. $\Delta \leftarrow \{v : v \notin \tilde{VC}, \deg(v) \geq d/2^i\}$;
4. $\tilde{VC} \leftarrow \tilde{VC} \cup \Delta$;
5. 从 G 中删除所有与 Δ 邻接的边;
6. **return** \tilde{VC} ;

- ▷ 运行时间： $d^{O(\log d)}$

- ▷ 近似比：($O(\log d), \epsilon n$)

209 / 250

问题5：Vertex Cover

想法：利用贪心算法设计 $(2, \epsilon n)$ 亚线性算法

- ▷ Vertex Cover v.s. Maximum Matching
- ▷ Vertex Cover v.s. Maximal Matching
- ▷ $MM \leq |VC| \leq 2|MM|$

贪心算法 (Greedy Sequential Matching Alg.)

1. $MM \leftarrow \emptyset$;
2. **for** 每一条边 $e \in E$ **do**
3. **if** u 和 v 都未匹配 **then**
4. 将边 e 的两个顶点加入 MM ;
5. **return** MM ;

212 / 250

问题5：Vertex Cover

- ▷ 如何利用贪心算法设计亚线性算法

Sublinear VC Alg II (Parnas-Ron framework)

1. 假设我们有一个 Oracle B ，给定 e 返回 $e \in MM$?
2. 独立均一地从 G 中选取 $s = \frac{8}{\epsilon^2}$ 个节点构成 S ;
3. **for** $v \in S$ **do**
4. **if** $\exists w \in N_v$ 使得 $B(v, w) = 1$ **then**
5. $X_v = 1$;
6. **else**
7. $X_v = 0$;
8. **return** $\hat{VC} = \frac{n}{2s} \cdot \sum_{v \in S} X_v + \frac{\epsilon}{2} n$;

Alg II 运行时间 $O(\text{Time}(B)/\epsilon^2)$

213 / 250

问题5：Vertex Cover

- ▷ 如何实现 Oracle B

Oracle B

1. 为查询边 e 随机生成一个 $\text{rank}_e \in [0, 1]$;
2. **for** 每条 e 的邻居边 $e' \in E$ **do**
3. **if** $\text{rank}_{e'} < \text{rank}_e$ **then**
4. 递归调用 Oracle B 检查 e' ;
5. **if** $e' \in MM$ **then return** 0;
6. **return** 1;

Oracle B 的期望运行时间 $2^{O(d)} \Rightarrow$ Alg II 运行时间 $\frac{2^{O(d)}}{\epsilon^2}$

215 / 250

- ▣ 概率基础
- ▣ 亚线性空间算法
- ▣ 亚线性时间算法
- ▣ 并行模型算法

237 / 250

并行计算模型

- 分布式计算环境（1k、1M台机器）如何计算？
- ▷ PRAM模型
 - ▷ Bulk Synch Parallel (BSP)模型
 - ▷ MapReduce模型

238 / 250

并行计算模型

MapReduce模型

- ▷ 所有数据形如 $\langle key, value \rangle$ ，计算分轮次进行
- ▷ 每一轮分为：Map、Shuffle、Reduce
- ▷ Map: 每个数据被 *map* 函数处理，输出一个新的数据集
- ▷ Shuffle: 对编程人员透明，所有在Map中输出的数据，被按照key分组，具备相同key的数据被分配到相同的reducer
- ▷ Reduce: 输入 $\langle k, v_1, v_2, \dots \rangle$ ，输出新的数据集

设计目标：更少的轮数、更少的内存、更少的工作量、更大的并行度

240 / 250

问题1：基本问题

构建倒排索引

- ▷ Map函数
输入： $\langle docID, content \rangle$
输出： $\langle word, docID \rangle$
- ▷ Reduce函数
输入： $\langle word, docID \rangle$
输出： $\langle word, list\ of\ docID \rangle$

241 / 250

问题1：基本问题

单词计数

- ▷ Map函数
输入： $\langle docID, content \rangle$
输出： $\langle word, 1 \rangle$
- ▷ Reduce函数
输入： $\langle word, 1 \rangle$
输出： $\langle word, count \rangle$

242 / 250

问题1：基本问题

检索Search

- ▷ Map函数
输入： $\langle (docID, lineno), content \rangle$
输出： $\langle docID, NULL \rangle$
- ▷ Reduce函数
输入： $\langle docID, (NULL, NULL, ...) \rangle$
输出： $\langle docID, NULL \rangle$

243 / 250

问题1：基本问题

矩阵乘法A

- ▷ Map()
- $((A, i, j), a_{ij}) \rightarrow (j, (A, i, a_{ij}))$
- $((B, j, k), b_{jk}) \rightarrow (j, (B, k, b_{jk}))$
- ▷ Reduce()
- $(j, (A, i, a_{ij}), (j, (B, k, b_{jk}))) \rightarrow ((i, k), a_{ij} * b_{jk})$
- ▷ Map+(): identity
- ▷ Reduce+()
- $((i, k), (v_1, v_2, \dots)) \rightarrow ((i, k), \sum v_i)$

244 / 250

问题1：基本问题

矩阵乘法B

- ▷ Map()
- $((A, i, j), a_{ij}) \rightarrow ((i, x), (A, j, a_{ij}))$ for all x
- $((B, j, k), b_{jk}) \rightarrow ((y, k), (B, j, b_{jk}))$ for all y
- ▷ Reduce()
- $((i, k), (A, j, a_{ij})) \wedge ((i, k), (B, j, b_{jk})) \rightarrow ((i, k), \sum a_{ij} * b_{jk})$

245 / 250

问题2：排序问题

使用 p 台处理器，输入 $\langle i, A[i] \rangle$

TeraSort: Round 1

- map: $\langle i, A[i] \rangle$
1. 输出 $\langle i \% p, ((i, A[i]), 0) \rangle$;
 2. 以概率 T/n ，这里 $T = \log(p)/\epsilon^2$
对所有 $j \in [0, p-1]$ ，输出 $\langle j, (A[i], 1) \rangle$;
- reduce: $\langle j, (A[i], y) \rangle$
1. 将 $y = 1$ 的数据收集为 S 并排序;
 2. 将 $y = 0$ 的数据收集为 B ;
 3. 构造 $(s_1, s_2, \dots, s_{p-1})$ ， s_j 为 S 中第 $j \lceil \frac{|S|}{p} \rceil$ ，
对 $(i, x) \in B$ 满足 $s_j < x \leq s_{j+1}$ ，输出 $\langle j, (i, x) \rangle$;

246 / 250

问题2：排序问题

使用 p 台处理器，输入 $\langle i, A[i] \rangle$

TeraSort: Round 2

- map: $\langle j, (i, x) \rangle$
1. 输出 $\langle j, (i, x) \rangle$;
- reduce: $\langle j, (i, x) \rangle$
1. 将所有 (i, x) 根据 x 排序并输出
- ▷ 随机算法
- ▷ $|S| \gg p$ ，样本点中的 $p-1$ 个分点
如果 s_j 是 S 中的 α 分点，那么以很高概率，它处于 $\alpha n \pm \epsilon n$

247 / 250