

Моменты и частые элементы

Всеволод Опарин

CS Клуб, Осень 2014

9 ноября

Задачи

- ▶ $\sigma = \langle (a_1, c_1), \dots, (a_m, c_m) \rangle$.
- ▶ $f_x = \sum_{i:a_i=x} c_i$ – частота.
- ▶ Кассовая модель: $c_i > 0$.
- ▶ Турникетная модель: c_i – произвольные.

Частота элемента

- ▶ По заданному x найти f_x .

k -ый момент

- ▶ Найти момент $F_k = \sum_{x:f_x>0} f_x^k$.

Эскиз (англ. sketch)

Определение

$DS(\sigma)$ – эскиз, если для двух потоков σ_1 и σ_2

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

COMB – эффективен по памяти.

Определение

Пусть $A \in \mathbb{R}^{k \times n}$, $\mathbf{f} \in \mathbb{R}^n$. $A\mathbf{f}$ – линейный эскиз.

Эскиз-счетчик

Идея

Берем случайную матрицу A , у которой в каждом столбце один ненулевой элемент: 1 или -1.

$$\begin{pmatrix} 0 & 0 & 0 & 0 & -1 & -1 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ -1 & 0 & 0 & 1 & 0 & 0 \end{pmatrix} \cdot \begin{pmatrix} f_1 \\ f_2 \\ \vdots \\ f_n \end{pmatrix} = \begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_k \end{pmatrix}$$

- ▶ Ответ для x : $\mathbf{v}^T \cdot A \cdot \mathbf{e}_x$.

Эскиз-счетчик

Алгоритм

- ▶ Пусть $k = \frac{1}{\epsilon^2}$.
- ▶ Выбираем две функции
 - ▶ $h : [n] \rightarrow [k]$,
 - ▶ $g : [n] \rightarrow \{-1, +1\}$из 2-независимых семейств.
- ▶ Заводим массив $v[1..k]$.
- ▶ Обработка пары (x, c) : $v[h(x)] += g(x) \cdot c$.
- ▶ Ответ для x : $g(x) \cdot v[h(x)]$.

Эскиз-счетчик

Анализ

- ▶ Спросили про элемент a .
- ▶ $Y_j = [h(a) = h(j)]$.
- ▶ $X = g(a) \sum_j f_j \cdot g(j) \cdot Y_j$.
- ▶ $\mathbf{E}[X] = f_a$.
- ▶ $\mathbf{D}[X] = \frac{\|\mathbf{f}\|_2^2 - f_a^2}{k}$.
- ▶ $\Pr[|X - f_a| \geq \varepsilon \cdot \|f_{-a}\|_2] \leq \frac{1}{3}$.

Повторить $O(\log \frac{1}{\delta})$ раз, взять медиану.

- ▶ Гарантия: $\Pr[|X - f_a| \geq \varepsilon \cdot \|f_{-a}\|_2] \leq \delta$.
- ▶ Память: $O(\frac{1}{\varepsilon^2} \log \frac{1}{\delta} \cdot (\log m + \log n))$.

Эскиз-мин. счетчик

Идея

Берем случайную матрицу A , у которой в каждом столбце ровно одна единица.

$$\begin{pmatrix} 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 \end{pmatrix} \cdot \begin{pmatrix} f_1 \\ f_2 \\ \vdots \\ f_n \end{pmatrix} = \begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_k \end{pmatrix}$$

- ▶ Ответ для x : $\mathbf{v}^T \cdot A \cdot \mathbf{e}_x$.

Эскиз-мин. счетчик

Алгоритм

- ▶ Пусть $k = \frac{1}{2 \cdot \epsilon}$.
- ▶ Выбираем одну функцию
 - ▶ $h : [n] \rightarrow [k]$,из 2-независимого семейства.
- ▶ Заводим массив $v[1..k]$.
- ▶ Обработка пары (x, c) : $v[h(x)] += c$.
- ▶ Ответ для x : $v[h(x)]$.

Эскиз-мин. счетчик

Анализ

- ▶ Спросили про элемент a .
- ▶ $Y_j = [h(a) = h(j)]$.
- ▶ $\text{ans} = f_a + X$.
- ▶ $X = \sum_{j \neq a} f_j \cdot Y_j$.
- ▶ $\mathbf{E}[X] = \frac{\|f_{-a}\|_1}{k}$.
- ▶ $\mathbf{Pr}[X \geq \varepsilon \cdot \|f_{-a}\|_1] \leq \frac{1}{2}$.

Повторить $O(\log \frac{1}{\delta})$ раз, взять минимум.

- ▶ Гарантия: $\mathbf{Pr}[|\text{ans} - f_a| \geq \varepsilon \cdot \|f_{-a}\|_1] \leq \delta$.
- ▶ Память: $O(\frac{1}{\varepsilon} \log \frac{1}{\delta} \cdot (\log m + \log n))$.

Момент F_2

Задача

- ▶ Дан поток $\sigma = \langle (a_1, c_1), \dots, (a_m, c_m) \rangle$. $f_x = \sum_{i:a_i=x} c_i$.
- ▶ Посчитать $F_2 = \sum f_x^2$.

Идея

Берем случайную матрицу $A \in \mathbb{R}^{t \times n}$, у которой все элементы из $\{-1, +1\}$.

$$\frac{1}{\sqrt{t}} \begin{pmatrix} 1 & -1 & -1 & -1 & 1 & 1 \\ -1 & 1 & -1 & 1 & -1 & -1 \\ 1 & -1 & 1 & 1 & 1 & -1 \\ 1 & 1 & 1 & 1 & -1 & 1 \\ 1 & -1 & -1 & 1 & 1 & -1 \end{pmatrix} \cdot \begin{pmatrix} f_1 \\ f_2 \\ \vdots \\ f_n \end{pmatrix} = \begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_t \end{pmatrix}$$

- ▶ Ответ: $\|\mathbf{v}\|_2^2$.

Момент F_2

Алгоритм

- ▶ Выбираем одну функцию
 - ▶ $g : [n] \rightarrow \{-1, +1\}$
из 4-независимого семейства.
- ▶ Заводим переменную $x = 0$.
- ▶ Обработка пары (x, c) : $x += c \cdot g(x)$.
- ▶ Ответ: x^2 .

Медиана средних

- ▶ $E[X] = Q$. Известна $D[X]$
- ▶ Цель – Z : $\Pr[|Z - Q| \geq \varepsilon \cdot Q] \leq \delta$.
- ▶ $Y = \frac{1}{k} \sum_{j=1}^k X_j$.
- ▶ $D[Y] = \frac{1}{k} D[X]$.
- ▶ $\Pr[|Y - Q| \geq \varepsilon \cdot Q] \leq \frac{D[X]}{k \cdot \varepsilon^2 Q^2}$.
- ▶ $k = \frac{3 \cdot D[X]}{Q^2 \varepsilon^2}$, чтобы получить $\frac{1}{3}$.
- ▶ Повторить $O(\log \frac{1}{\delta})$ раз, взять медиану.
- ▶ Ухудшение по памяти: $O(\varepsilon^{-2} \log \frac{1}{\delta} \frac{D[X]}{Q^2})$.

Момент F_2

Анализ

- ▶ $Y_j = h(j)$
- ▶ $X = \sum_j f_j \cdot Y_j$.
- ▶ $\mathbf{E} [X^2] = F_2$.
- ▶ $\mathbf{E} [X^4] = F_4 + 6 \cdot \sum_{i < j} f_i^2 \cdot f_j^2$.
- ▶ $\mathbf{D} [X^2] = 2 \cdot F_2^2$.

Применяем медиану средних.

- ▶ Гарантия: $\Pr [\left| \|\mathbf{v}\|_2^2 - F_2 \right| \geq \varepsilon \cdot F_2] \leq \delta$.
- ▶ Память: $O\left(\frac{1}{\varepsilon^2} \log \frac{1}{\delta} (\log m + \log n)\right)$.

Регрессия

Задача

- ▶ Y – зависимая переменная.
- ▶ A_j – предикторы.
- ▶ Хотим построить линейную зависимость
$$Y = x_1 \cdot A_1 + \dots + x_d \cdot A_d.$$
- ▶ Наблюдение за предикторами дается матрицей A .
Наблюдения за Y – вектором y .
- ▶ Найти $x^* = \arg \min \|Ax - y\|_2$.

Нужно решить $A^T Ax = A^T y$, а наблюдений много...

Регрессия

Идея

- ▶ Пусть $m = O(d\varepsilon^{-1} \log \frac{1}{\delta})$.
- ▶ Возьмем матрицу $S \in \{-1, +1\}^{n \times m}$.
- ▶ Матрица должна быть $\Omega(d + \log \frac{1}{\delta})$ -независимой.
- ▶ Найдем $\hat{x} = \arg \min \|S^T(Ax - y)\|_2$.
- ▶ Гарантия: $\Pr [\|A\hat{x} - y\|_2 \leq (1 + \varepsilon)\|Ax^* - y\|_2] \geq 1 - \delta$.
- ▶ Память: $O(d^2\varepsilon^{-1} \log \frac{1}{\delta} \log N)$.

Момент F_k

Задача

- ▶ Дан поток $\sigma = \langle a_1, \dots, a_m \rangle$, $a_i \in [n]$.
- ▶ Найти $F_k = \sum_{x: f_x > 0} f_x^k$.

Алгоритм

- ▶ Взять случайный элемент a в потоке.
- ▶ Посчитать, сколько раз он встретился. Счетчик в r .
- ▶ Вернуть $m \cdot (r^k - (r - 1)^k)$.

Момент F_k

Анализ

- ▶ A – выбранный элемент. R – значение r в конце.
- ▶ Выбираем элемент, выбираем появление.
- ▶ $\mathbf{E}[X] = F_k$.
- ▶ $\mathbf{D}[X] \leq \mathbf{E}[X^2] = k \cdot F_1 \cdot F_{2k-1}$.

Лемма

Для $n > 0$, $x_i \geq 0$, $k \geq 1$

$$\left(\sum x_i\right)\left(\sum x_i^{2k-1}\right) \leq n^{1-\frac{1}{k}}\left(\sum x_i^k\right)^2.$$

Суммы берутся по $i \in [n]$.

- ▶ $\mathbf{D}[X] \leq kn^{1-\frac{1}{k}}F_k^2$.
- ▶ Память: $O(\varepsilon^{-2} \cdot \log \frac{1}{\delta} \cdot k \cdot n^{1-\frac{1}{k}}(\log m + \log n))$

Спасибо!

Вопросы?