

Sublinear Algorithms for Big Data

Exam Problems

Grigory Yaroslavtsev

<http://grigory.us>



Problem 1: Modified Chernoff Bound

- **Problem:** Derive Chernoff Bound #2 from Chernoff Bound #1. Explain your answer in detail.
- **(Chernoff Bound #1)** Let $X_1 \dots X_t$ be **independent and identically distributed** r.v.s with range $[0, 1]$ and expectation μ . Then if $X = \frac{1}{t} \sum_i X_i$ and $1 > \delta > 0$,

$$\Pr[|X - \mu| \geq \delta\mu] \leq 2 \exp\left(-\frac{t\mu\delta^2}{3}\right)$$

- **(Chernoff Bound #2)** Let $X_1 \dots X_t$ be **independent and identically distributed** r.v.s with range $[0, c]$ and expectation μ . Then if $X = \frac{1}{t} \sum_i X_i$ and $1 > \delta > 0$,

$$\Pr[|X - \mu| \geq \delta\mu] \leq 2 \exp\left(-\frac{t\mu\delta^2}{3c}\right)$$

Problem 2: Modified Chebyshev

- **Problem:** Derive Chebyshev Inequality #2 from Chebyshev Inequality #1. Explain your answer in detail.

- **(Chebyshev Inequality #1)** For every $c > 0$:

$$\Pr \left[|\mathbf{X} - \mathbb{E}[\mathbf{X}]| \geq c \sqrt{\text{Var}[\mathbf{X}]} \right] \leq \frac{1}{c^2}$$

- **(Chebyshev Inequality #2)** For every $c' > 0$:

$$\Pr \left[|\mathbf{X} - \mathbb{E}[\mathbf{X}]| \geq c' \mathbb{E}[\mathbf{X}] \right] \leq \frac{\text{Var}[\mathbf{X}]}{(c' \mathbb{E}[\mathbf{X}])^2}$$

Problem 3: Sparse Recovery Error

- **Definition (frequency vector):** f_i = frequency of i in the stream = # of occurrences of value i

$$f = \langle f_1, \dots, f_n \rangle$$

- **Sparse Recovery:** Find g such that $\|f - g\|_1$ is minimized among g 's with at most k non-zero entries.
- **Definition:** $Err^k(f) = \min_{g: \|g\|_0 \leq k} \|f - g\|_1$
- **Problem:** Show that $Err^k(f) = \sum_{i \notin S} |f_i|$ where S are indices of k largest f_i . Explain your answer **formally** and **in detail**.

Problem 4: Approximate Median Value

- Stream: m elements x_1, \dots, x_m from universe $[n] = \{1, 2, \dots, n\}$.
- $S = \{x_1, \dots, x_m\}$ (all distinct) and let
$$\text{rank}(y) = |\{x \in S : x \leq y\}|$$
- Median $M = x_i$, where $\text{rank}(x_i) = \frac{m}{2} + 1$ (m odd).
- **Algorithm:** Return y = the median of a sample of size t taken from S (with replacement).
- **Problem:** Does this algorithm give a 10% approximate value of the median with probability $\geq \frac{2}{3}$, i.e. y such that

$$M - \frac{n}{10} < y < M + \frac{n}{10}$$

if $t = o(n)$? Explain your answer.

Problem 5: Lower bound on F_k

- **Definition (frequency vector):** f_i = frequency of i in the stream = # of occurrences of value i

$$f = \langle f_1, \dots, f_n \rangle$$

- **Define (frequency moment):** $F_k = \sum_i f_i^k$
- **Problem:** Show that for all integer $k \geq 1$ it holds that:

$$F_k \geq n \left(\frac{m}{n} \right)^k$$

- **Hint:** worst-case when $f_1 = \dots = f_n = \frac{m}{n}$. Use convexity of $g(x) = x^k$

Problem 6: Approximate MST Weight

- Let G be a weighted graph with weights of all edges being integers between 1 and W .
- **Definition:** Let n_i be the # of connected components if we remove all edges with weight $> (1 + \epsilon)^i$.
- **Problem:** For some constant $C > 0$ show the following bounds on the weight of the minimum spanning tree $w(MST)$:

$$\begin{aligned} w(MST) &\leq \sum_{i=0}^{\lceil \log_{1+\epsilon} W \rceil} \epsilon (1 + \epsilon)^i n_i \\ &\leq (1 + C \epsilon) w(MST) \end{aligned}$$

Evaluation

- Deadline: September 1st , 2014, 23:59 GMT
- Submissions in English: grigory@grigory.us
 - Submission Title (once): Exam + Space + “Your Name”
 - Question Title: Question + Space + “Your Name”
 - Submission format
 - PDF from LaTeX (best)
 - PDF
- You can work in groups (up to 3 people)
 - Each group member lists all others on the submission
- Point system
 - Course Credit = 2 points
 - Problems 1-3 = 0.5 point each
 - Problem 4 = 1 point
 - Problems 5-6 = 2 point each