

Sublinear Algorithms for Big Data

Lecture 3

Grigory Yaroslavtsev

<http://grigory.us>



SOFSEM 2015



- URL: <http://www.sofsem.cz>
- 41st International Conference on Current Trends in Theory and Practice of Computer Science (SOFSEM'15)
- When and where?
 - January 24-29, 2015. Czech Republic, Pec pod Snezkou
- Deadlines:
 - August 1st (tomorrow!): Abstract submission
 - August 15th: Full papers (proceedings in LNCS)
- I am on the Program Committee ;)

Today

- Count-Min (continued), Count Sketch
- Sampling methods
 - ℓ_2 -sampling
 - ℓ_0 -sampling
- Sparse recovery
 - ℓ_1 -sparse recovery
- Graph sketching

Recap

- Stream: elements from universe $[]$
= $\{1, 2, \dots\}$, e.g.
 $\langle x_1, x_2, \dots \rangle = \langle 5, 8, 1, 1, 1, 4, 3, 5, \dots, 10 \rangle$
- f_x = frequency of x in the stream = # of occurrences of value x , $f = \langle f_1, f_2, \dots \rangle$

Count-Min

- $h_1, \dots, h_k : [n] \rightarrow [m]$ are 2-wise independent hash functions
- Maintain m counters with values:
 - $C_j = \#$ elements in the stream with $h_i(x) = j$
- For every x the value $C_{h_i(x)} \geq f(x)$ and so:
 - $\tilde{f} = \min_{1 \leq i \leq k} C_{h_i(x)}$
- If $m = \frac{2}{\epsilon}$ and $k = \log_2 \frac{1}{\epsilon}$ then:

$$| \tilde{f} - f(x) | \leq \epsilon$$

More about Count-Min

- **Authors:** Graham Cormode, S. Muthukrishnan [LATIN'04]
- Count-Min is linear:
$$\text{Count-Min}(S1 + S2) = \text{Count-Min}(S1) + \text{Count-Min}(S2)$$
- Deterministic version: CR-Precis
- Count-Min vs. Bloom filters
 - Allows to approximate values, not just 0/1 (set membership)
 - Doesn't require mutual independence (only 2-wise)
- FAQ and Applications:
 - <https://sites.google.com/site/countminsketch/home/>
 - <https://sites.google.com/site/countminsketch/home/faq>

Fully Dynamic Streams

- Stream: updates $(i, \Delta) \in [1, m] \times \mathbb{R}$ that define vector \mathbf{v} where $v_i = \sum \Delta$.

- **Example:** For $m = 4$

$$\langle (1,3), (3, 0.5), (1,2), (2, -2), (2,1), (1, -1), (4,1) \rangle \\ = (4, -1, 0.5, 1)$$

- Count Sketch: Count-Min with **random signs** and **median** instead of min:

Count Sketch

- In addition to $:\mathbb{R}^d \rightarrow \mathbb{R}^k$ use random signs $[\] \rightarrow \{-1, 1\}$

$$y_i = \sum_{j=1}^d s_{ij} x_j$$

- Estimate:

$$\hat{x}_i = \frac{1}{k} \sum_{j=1}^k y_j s_{ij}$$

- Parameters: $k = \left(\log \frac{1}{\epsilon} \right), \quad m = \frac{3}{\epsilon}$

$$\Pr \left[\left| \frac{\hat{x}_i}{x_i} - 1 \right| \geq \epsilon \right] \leq \epsilon$$

ℓ -Sampling

- Stream: updates $(i, \Delta) \in [n] \times \mathbb{R}$ that define vector x where $x_i = \sum_{(i, \Delta) \in \text{updates}} \Delta$.

- ℓ -Sampling: Return random $i \in [n]$ and $w \in \mathbb{R}$:

$$\Pr [|w - x_i| \leq \epsilon] = (1 \pm \epsilon) \frac{|x_i|}{\sum_{i \in [n]} |x_i|} \pm \epsilon$$

$$= (1 \pm \epsilon)$$

Application: Social Networks

- Each of n people in a social network is friends with some arbitrary set of other $n - 1$ people
- Each person knows only about their friends
- With no communication in the network, each person sends a postcard to Mark Z.
- If Mark wants to know if the graph is connected, how long should the postcards be?

Optimal estimation

- Yesterday: (ϵ, δ) -approximate

- $\tilde{O}\left(\frac{1}{\epsilon^2} \log \frac{1}{\delta}\right)$ space for $\sum_{i=1}^n |x_i|^2$

- $\tilde{O}\left(\log \frac{1}{\delta}\right)$ space for σ^2

- **New algorithm:** Let (x, y) be an ℓ_2 -sample.

Return $\hat{\sigma}^2 = \hat{\sigma}_2^2 - 2$, where $\hat{\sigma}_2^2$ is an ϵ -estimate of σ^2

- **Expectation:**

Optimal estimation

- **New algorithm:** Let (x, y) be an ℓ_2 -sample.

Return $\hat{\sigma}_2^2 = \hat{\sigma}_2^{-2}$, where $\hat{\sigma}_2$ is an ± 2 estimate of σ_2

- **Variance:**

$$\mathbb{E}[\hat{\sigma}_2^2] \leq \mathbb{E}[\sigma_2^2] = \sum_{x \in [1, 2]} \mathbb{E}[\sigma_2^2 | x]$$

$$= \pm 2 \sum_{x \in [1, 2]} \frac{1}{2} \sigma_2^2 (1 - 2^{-x}) = \pm 2 \frac{1}{2} \sigma_2^2 (1 - 2^{-2}) \leq \pm 2 \sigma_2^2$$

- **Exercise:** Show that $\mathbb{E}[\sigma_2^2] \leq 1 - \frac{2}{2} \sigma_2^2$

ℓ_2 -Sampling: Basic Overview

- Assume $\sum_{i=1}^n w_i^2 = 1$. Weight w_i by $\sqrt{\frac{1}{w_i}}$, where $w_i \in [0, 1]$:

$$= \left(\frac{1}{w_1} \frac{1}{w_2} \dots \frac{1}{w_n} \right)$$

$$= \left(\frac{1}{w_1} \frac{1}{w_2} \dots \frac{1}{w_n} \right)$$
 where $\frac{1}{w_i} = \sqrt{\frac{1}{w_i}}$
- For some value ϵ , return (i, j) if there is a unique i such that $w_i^2 \geq \epsilon$
- Probability (i, j) is returned if w_i is large enough:

ℓ_2 -Sampling: Part 1

- Use Count-Sketch with parameters (k, b) to sketch
- To estimate $\|x\|_2^2$:

$$\|x\|_2^2 = \sum_{i=1}^n x_i^2 \quad \text{and} \quad \hat{\|x\|_2^2} = \frac{1}{k} \sum_{j=1}^b \sum_{i=1}^n x_i^2 \mathbb{1}_{h(i)=j}$$

- **Lemma:** With high probability if $k = \Theta(\log n)$

$$\hat{\|x\|_2^2} = \|x\|_2^2 \pm \epsilon \|x\|_2^2 \pm \frac{\epsilon}{k} \sum_{i=1}^n x_i^2$$
- **Corollary:** With high probability if $k = \Theta(\log n)$ and $n \gg \frac{1}{\epsilon^2}$,

Proof of Lemma

- Let $\epsilon = \frac{1}{3}$
- By the analysis of Count Sketch $\left[\sum_{i=1}^n x_i^2 \right] \leq \frac{2}{\epsilon^2} \sum_{i=1}^n x_i$ and by Markov:

$$\Pr \left[\sum_{i=1}^n x_i^2 \geq \frac{3}{\epsilon^2} \sum_{i=1}^n x_i \right] \leq \frac{\epsilon^2}{3}$$
- If $\left| \sum_{i=1}^n x_i \right| \geq \frac{2}{\epsilon} \sum_{i=1}^n x_i$, then $\left| \sum_{i=1}^n x_i \right|^2 = \frac{4}{\epsilon^2} \left(\sum_{i=1}^n x_i \right)^2$
- If $\left| \sum_{i=1}^n x_i \right| \leq \frac{2}{\epsilon} \sum_{i=1}^n x_i$, then

ℓ_2 -Sampling: Part 2

- Let $x_i = 1$ if $\hat{x}_i^2 \geq \frac{4}{n}$ and $x_i = 0$ otherwise
- If there is a unique i with $x_i = 1$ then return (i, \hat{x}_i^2) .
- Note that if $\hat{x}_i^2 \geq \frac{4}{n}$ then $\frac{1}{n} \leq \frac{\hat{x}_i^2}{4}$ and so

$$\hat{x}_i^2 = x_i^2 \pm \frac{1}{n} = x_i^2 \pm \frac{\hat{x}_i^2}{4},$$

therefore $x_i^2 = \pm 4 \hat{x}_i^2$

Proof of Lemma

- Let $t = \frac{4}{\dots}$. We can upper-bound $\Pr[\dots = 1]$:

$$\Pr[\dots = 1] = \Pr\left[\sum_{i=1}^n \dots^2 \geq \dots\right]$$

$$\leq \Pr\left[\frac{\dots}{-4 \dots} \geq \dots\right] \leq \dots$$

Similarly, $\Pr[\dots = 1] \geq \dots$.

- Using independence of \dots , probability of unique \dots with $\dots = 1$:

$$\sum \Pr[\dots = 1 \dots = 0] \geq \sum \Pr[\dots = 1] \left(1 - \sum \Pr[\dots = 1]\right)$$

Proof of Lemma

- Let $t = \frac{4}{\dots}$. We can upper-bound $\Pr \left[\dots = 1 \right]$:

$$\Pr \left[\dots = 1 \right] = \Pr \left[\dots^2 \geq \dots \right]$$

$$\leq \Pr \left[\frac{\dots}{-4} \geq \dots \right] \leq \dots$$

Similarly, $\Pr \left[\dots = 1 \right] \geq \dots$.

- We just showed:

$$\sum \Pr \left[\dots = 1, \dots = 0 \right] \approx 1 /$$

ℓ_0 -sampling

- Maintain $\tilde{\sigma}$ and (1 ± 0.1) -approximation to σ
- Hash items using $h : [] \rightarrow [0, 2^{-1}]$ for $\in [\log]$
- For each σ , maintain:

$$= (1 \pm 0.1) |\{ \sigma \mid h(\sigma) = 0 \}|$$

$$= \sum_{\sigma, h(\sigma) = 0}$$

$$= \sum_{\sigma, h(\sigma) = 0}$$

Proof of Lemma

- Let $n = \lceil \log_2 \tilde{n} \rceil$ and note that $2^{n-1} < 2^n < 12 \cdot 2^{n-1}$
- For any i , $\Pr [h(i) = 0] = \frac{1}{2}$
- Probability there exists a unique i such that $h(i) = 0$,

$$\sum \Pr [h(i) = 0 \quad \forall i \neq j, h(i) \neq 0]$$

$$= \sum \Pr [h(i) = 0] \Pr [\forall i \neq j, h(i) \neq 0 \mid h(i) = 0]$$

Sparse Recovery

- **Goal:** Find \hat{g} such that $\|y - A\hat{g}\|_1$ is minimized among g 's with at most k non-zero entries.
- **Definition:** $\hat{g} = \min_{g: \|g\|_0 \leq k} \|y - Ag\|_1$
- **Exercise:** $\hat{g} = \sum_{i \in S} |g_i| e_i$ where S are indices of largest k entries of $|Ag|$
- Using $\left(\frac{1}{k} \log \binom{n}{k} \right)$ space we can find \hat{g} such

Count-Min Revisited

- Use Count-Min with $\epsilon = \frac{1}{4}$, $d = 4$
- For $i \in [n]$, let $\tilde{x}_i = \sum_{j \in [d]} h_j(x_i)$ for some row $j \in [d]$
- Let $S = \{i_1, \dots, i_k\}$ be the indices with max. frequencies. Let E be the event there doesn't exist $i \in S$ with $h_j(i) = h_j(i)$

- Then for $i \in [n]$:

$$\Pr \left[\left| \tilde{x}_i - x_i \right| \geq \frac{\epsilon}{4} \right] =$$

$$\Pr \left[\text{not } E \right] \times \Pr \left[\left| \tilde{x}_i - x_i \right| \geq \frac{\epsilon}{4} \mid \text{not } E \right]$$

+

Sparse Recovery Algorithm

- Use Count-Min with $\epsilon = \frac{1}{\log n}$, $m = 4 / \epsilon^2$
- Let $\tilde{f} = (\tilde{f}_1, \tilde{f}_2, \dots, \tilde{f}_n)$ be frequency estimates:

$$| \tilde{f}_i - f_i | \leq \frac{\epsilon}{k}$$
- Let \tilde{f}^k be \tilde{f} with all but the k -th largest entries replaced by 0.
- **Lemma:** $\| \tilde{f}^k - f \|_1 \leq (1 + 3 \frac{\epsilon}{k}) \| f \|_1$

$$\| \tilde{\cdot} - \cdot \|_1 \leq (1 + 3) \quad ()$$

- Let $\cdot, \tilde{\cdot} \subseteq []$ be indices corresponding to largest value of and $\tilde{\cdot}$.

- For a vector $\in \mathbb{R}$ and $\subseteq []$ denote as the vector formed by zeroing out all entries of except for those in \cdot .

$$\begin{aligned} \| \cdot - \tilde{\cdot} \|_1 &\leq \| \cdot - \cdot \|_1 + \| \tilde{\cdot} - \cdot \|_1 \\ &= \| \cdot \|_1 + \| \tilde{\cdot} - \cdot \|_1 \\ &= \| \cdot \|_1 + \| \tilde{\cdot} - \cdot \|_1 + \left(\| \cdot \|_1 + \| \tilde{\cdot} - \cdot \|_1 \right) \\ &\leq \| \cdot \|_1 + \| \tilde{\cdot} - \cdot \|_1 + 2 \| \cdot \|_1 + \| \tilde{\cdot} - \cdot \|_1 \end{aligned}$$

Thank you!

- Questions?