# Lecture 8: Streaming and Sketching I

*"In low-space, nobody can remember your stream."*

> We will follow for this chapter the (excellent) lecture notes by Amit Chakrabarti [AC], available at `https://www.cs.dartmouth.edu/~ac/Teach/data-streams-lecnotes.pdf`.

## The Basic Setup

We will specifically focus on one-pass algorithms, unless specified otherwise. $m$ denotes the length of the stream

$$\sigma = \langle a_1, \ldots, a_m \rangle$$

where each $a_i$ belongs to the universe $\mathcal{X}$ of size $n$. We do not impose any bound on the time complexity of our algorithms, but we will enforce that they use very little memory (space), with space complexity denoted by $s$. We will aim for

$$s = o(\min(m, n))$$

and would love to use much less, ideally

$$s = O(\log m + \log n)$$

or, if not, $s = \text{poly}(\log m, \log n)$. To do so, we will allow for randomised algorithms *and* approximation algorithms, where the quality of the approximation will be controlled by a parameter $\varepsilon > 0$, usually thought of as an (arbitrarily) small fixed constant.

## The Majority Problem

Chapter 1 of [AC]

**Theorem 39.** *The* MISRA-GRIES *algorithm is a* deterministic *one-pass algorithm which, for any given parameter $\varepsilon \in (0, 1]$, provides $\hat{f}_1, \ldots, \hat{f}_n$ of all element frequencies such that*

$$f_j - \varepsilon m \le \hat{f}_j \le f_j, \qquad j \in [n]$$

*with space complexity $s = O(\log(mn)/\varepsilon)$. (In particular, it can be used to solve the* MAJORITY *problem in two passes.)*

## The Approximate Counting Problem

We will describe and analyse the Morris Counter algorithm, due to, well, Morris, which provides a constant-factor estimate of the number of elements of the stream: that is, an $F_1$ estimator.

**Theorem 40.** *The medians-of-means version of the* MORRIS COUNTER *is a* randomised *one-pass algorithm which, for any given parameters* $\varepsilon, \delta \in (0,1]$, *provides an estimate* $\widehat{d}$ *of the number* $d$ *of non-zero elements of the stream such that*

$$\Pr\Big[ (1-\varepsilon)d \le \widehat{d} \le (1+\varepsilon)d \Big] \ge 1 - \delta$$

*with space complexity*

$$s = O\left( \frac{\log \log m}{\varepsilon^2} \cdot \log \frac{1}{\delta} \right)$$

*that is,* doubly logarithmic *in* $m$.

But we can do better!

**Theorem 41.** *The "careful" version of* MORRIS COUNTER *is a* randomised *one-pass algorithm which, for any given parameters* $\varepsilon, \delta \in (0,1]$, *provides an estimate* $\widehat{d}$ *of the number* $d$ *of non-zero elements of the stream such that*

$$\Pr\Big[ (1-\varepsilon)d \le \widehat{d} \le (1+\varepsilon)d \Big] \ge 1 - \delta$$

*with space complexity*

$$s = O\left( \log \log m + \log \frac{1}{\varepsilon} + \log \frac{1}{\delta} \right)$$

*that is,* doubly logarithmic *in* $m$ *and* logarithmic *in* $1/\varepsilon$.

## The Distinct Elements Problem

We start this section with the TIDEMARK algorithm, due to Alon, Matias and Szegedy (AMS), which provides a constant-factor estimate of the number of distinct elements of the stream: that is, an $F_0$ estimator.

1: Pick $h\colon [n] \to [n]$ from a strongly universal hashing family
2: $z \leftarrow 0$
3: **for all** $1 \le i \le m$ **do**
4:     Get item $a_i \in [n]$
5:     **if** $\operatorname{zeros}(h(a_i)) \ge z$ **then**
6:         $z \leftarrow \operatorname{zeros}(h(a_i))$
7: **return** $\sqrt{2} \cdot 2^z$

**Theorem 42.** *The (median trick version of the)* TIDEMARK *(AMS) algorithm is a* randomised *one-pass algorithm which, for any given*

*parameter $\delta \in (0, 1]$, provides an estimate $\widehat{d}$ of the number $d$ of distinct elements of the stream such that, for some absolute constant $C > 0$,*

$$\Pr\left[\frac{1}{C} \cdot d \leq \widehat{d} \leq C \cdot d\right] \geq 1 - \delta$$

*with space complexity*

$$s = O\left(\log n \cdot \log \frac{1}{\delta}\right).$$

This is not bad, but can we achieve estimation factor arbitrarily close to one, say, $1 + \varepsilon$? The answer is yes: the following algorithm, due to Bar-Yossef, Jayram, Kumar, Sivakumar and Trevisan (BJKST), does exactly that.

Algorithm 16: The BJKST algorithm

---

**Input:** Parameter $\varepsilon \in (0, 1]$
1: Set $k \leftarrow O(\log^2 n / \varepsilon^4)$, $T \leftarrow \Theta(1/\varepsilon^2)$
2: Pick $h : [n] \rightarrow [n]$ from a strongly universal hashing family
3: Pick $g : [n] \rightarrow [k]$ from a strongly universal hashing family

4: $z \leftarrow 0$, $B \leftarrow \emptyset$
5: **for all** $1 \leq i \leq m$ **do**
6:     Get item $a_i \in [n]$
7:     **if** $\text{zeros}(h(a_i)) \geq z$ **then**
8:         $B \leftarrow B \cup \{(g(a_i), \text{zeros}(h(a_i)))\}$
9:         **while** $|B| \geq T$ **do**
10:             $z \leftarrow z + 1$
11:             Remove every $(a, b)$ with $b < z$ from $B$
12: **return** $|B| \cdot 2^z$

---

**Theorem 43.** *The (median trick version of the)* BJKST *algorithm is a randomised one-pass algorithm which, for any given parameters $\varepsilon, \delta \in (0, 1]$, provides an estimate $\widehat{d}$ of the number $d$ of distinct elements of the stream such that, for some absolute constant $C > 0$,*

$$\Pr\left[(1 - \varepsilon) \cdot d \leq \widehat{d} \leq (1 + \varepsilon)d\right] \geq 1 - \delta$$

*with space complexity*

$$s = O\left(\left(\log n + \frac{\log(1/\varepsilon) + \log\log n}{\varepsilon^2}\right) \cdot \log\frac{1}{\delta}\right).$$

This is pretty good, but... *Is it optimal?*