

Today:

- Nice properties of linear sketching algorithms
- Streaming algorithms
- Frequency moments F_p
- AMS sketch for F_2

Reminder: Linear sketching algorithms
(such as Count-Min Sketch)

Can be seen as:

$$\left[\begin{array}{c} \end{array} \right] = \left[\begin{array}{c} \text{randomized} \\ \text{matrix} \end{array} \right] \left[\begin{array}{c} \end{array} \right]$$

Diagram illustrating the components of a linear sketching algorithm:

- The leftmost bracket represents the "maintained low-dimensional sketch".
- The middle bracket represents the "what the algorithm does".
- The rightmost bracket represents the "high dimensional data / input".
- A red arrow points from the "what the algorithm does" bracket to the "randomized matrix" in the equation.
- A red arrow points from the "high dimensional data / input" bracket to the rightmost bracket in the equation.

In CMS:

$$\left[\begin{array}{c} \end{array} \right] = \left[\begin{array}{c} \text{randomized} \\ \text{matrix} \end{array} \right] \left[\begin{array}{c} \end{array} \right]$$

↑ ↑
kt counts 0/1 - matrix
Corresponding describing how
to buckets values get assigned
 to buckets
 by hash functions frequency
 counts f(· · ·)
 interpreted
 as a vector

Nice properties of linear sketching algorithms:

(follow from basic linear algebra properties)

- can handle both insertions & deletions

$$s = M d \leftarrow \text{data}$$

↑ ↑
sketch randomized
 matrix

updated sketch change in the vector
↓ corresponding to data insertion

Insertion handling: $s' = M(d + \Delta) = Md + M\Delta = s + M\Delta$

Deletion handling: compute $M\Delta$ but subtract from sketch

- can handle multiples of updates easily:

Instead adding $M\Delta$ to the sketch add $\alpha M\Delta$, where α is how many times the operation has to be applied.

Example: adding 1,000,000 copies of x in CMS takes essentially the same amount of time as adding one copy of x

- can handle distributed data:

$$d = d_1 + d_2 + \dots + d_k$$

each part in a different location

- no need to send all data to a single location
- each site computes $s_i \stackrel{\text{def}}{=} M d_i$ and sends it to a central location
- central location can now add sketches to get sketch

$$\begin{aligned}
 \text{for all data: } S &\stackrel{\text{def}}{=} \sum s_i = \\
 &= \sum M d_i \\
 &= M \sum d_i \\
 &= M d
 \end{aligned}$$

Important concept: Streaming algorithms

Streaming algorithm

Big stream of data
 $a_1, a_2, a_3, \dots, a_n$
 individual items,
 e.g., database records,
 graph edges, numbers, ...

- The algorithm reads and processes the input stream items one by one
- The algorithm should use much less space than the input size
- Natural question: How much space needed to solve a specific problem?

Comments:

- CMS can be seen as a streaming algorithm
- Insertion-only vs. insertions & deletions:

In the simplest version, items just arrive, but in the more general version (including CMS), each input item is either "insert x " or "delete x " for some x .

- One pass vs. multiple passes:

As presented, the algorithm gets to read each element once, but in some scenarios, it makes sense to assume that multiple passes are allowed. Example: If the data is stored on an external storage device, then sequential reading may maximize throughput (i.e.,

the rate at which data is read and processed). In this case, perhaps a small number of passes over the data is achievable.

Next problem: Frequency moments

Setting (same as in CMS for simplicity):

Our data is a set of items in X .

$f(x)$, for $x \in X$, is the number of copies of x .

p -th moment:

$$F_p(\text{our data}) = \sum_{x \in X} |f(x)|^p$$

$\nearrow p \in (0, \infty)$

\nearrow absolute value

because in some settings it makes sense to consider negative values

In the limit, if we assume that $0^0 = 0$ (and $x^0 = 1$ for $x \neq 0$), then F_0 is the number of distinct elements in our set.

Plan for next lectures:

- Approximating F_2 (AMS sketch)
- Approximating $F_0 = \# \text{ distinct elements}$