# Lecture 9

## Streaming algorithms

Also good introduction to sampling
and estimation algorithms.

## Setting

a set of objects/items/tokens
come in a stream/online fashion

$$e_1, e_2, \ldots, e_m$$

## Examples:
- $e_i$ is a number from $[n]$
- $e_i$ is an edge $(u, v)$ in a graph
- $e_i$ is a vector/point in $R^d$
- $e_i$ is a row/column of a matrix

- $e_i$ is a matrix.

Assumption: $m$ is large and unknown. Cannot afford to store all of $e_1, \ldots, e_m$ in memory.

We have say, $B$ space where $B \ll m$ (assume one token is one unit for now).

What functions of the input can we compute?

Depends on $B$ and we have various tradeoff between efficiency, accuracy etc.

# Simple but powerful setting.

each $e_i \in [n]$ hence a non-neg
integer from say large range $0..n-1$
or $1..n$.

Ex: $n = 1000$

## stream

$$5, 3, 1, 1, 2, 2, 10, 5, 90$$

Implicitly defines a vector $\bar{f} \in \mathbb{Z}^n$
that starts at all $\bar{0}$ vector

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Each $e_i$ adds 1 to vector $\bar{f}_{e_i}$

$$5, 3, 1, 1, 2, 2, 10, 5, 90$$

$$\begin{bmatrix} 0 \\ 0 \\ . \\ . \\ . \\ 0 \end{bmatrix} \xrightarrow{\quad} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ . \\ . \\ . \end{bmatrix} \xrightarrow{\quad} \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 1 \\ . \\ . \\ . \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ . \\ \\ \\ \end{bmatrix} \quad -$$

## Frequency moment estimation

Given stream $it$ define $\bar{f}$ at the end of the stream.

Want to estimate / compute $k$ th
frequency moment

$$F_k = \sum_{i=1}^{n} f_i^k \quad \text{or} \quad \left( \sum_{i=1}^{n} f_i^k \right)^{\frac{1}{k}} .$$

- $k = 0$    # of distinct elements

- $k = 1$    m.    <u>easy</u>

- $k = 2$    $l_2$ norm which is very important as we will see.

- $k = \infty$    largest coordinate / heavy hitters.

$k = 0, \ k = 2, \ k = \infty$ most important.

If we can store all $c_i$ or $f$ then problems are of course

trivial.

**Question** can we compute/estimate
with $B \ll m$.

Yes! with $B = \tilde{O}(1)$!
But requires _randomization_ and
_approximation_. Can show that
without either not feasible
with $o(n)$ memory.

# Reservoir Sampling

Given stream $e_1, e_2, \ldots, e_m$ with $m$ unknown. Want to have a uniformly random sample.

- $S \leftarrow \emptyset$
- $m \leftarrow 0$
- While (stream has not ended)

  $$m \leftarrow m+1$$
  $$e_m \leftarrow \text{current element}$$
  with prob $\frac{1}{m}$.
  $$S \leftarrow e_m$$

  end while
- output $S$

If we want $k$ samples with replacement can use above independently.

$k$ samples <u>without</u> replacement.

$$S \leftarrow \phi, \quad m \leftarrow \phi$$

While (stream is not done)

  $m \leftarrow m+1$

  if $m \leq k$ add $e_m$ to $S$.

  else

    with prob $\frac{k}{m}$

      replace a random element of $S$ with $e_m$

end while

Output $S$.

Exercise: Prove correctness.

See notes to learn about "weighted" sampling.

# Distinct Element Estimation

Given stream $e_1, e_2, \ldots, e_m$ where each $e_i \in [n]$, want to know $F_0$, # of distinct elements seen.

Ex: 1, 1, 3, 1, 1, 1, 1, 1, 3, 5, 5, 5, 1, 2

4 is # of distinct elements.

$n$ is large $m$ is large.

Ex: An interesting application.

Packets through a high-speed internet switch.

(src, dest, content)

$\nearrow$ IP address of source and destination

128 bits in IPv6.

so $\overline{\overline{n}}$ is $2^{128}$ technically.

Q: How many distinct source address?

## Offline Solution

Can store distinct elements using
a dictionary data structure.
$\Theta(k)$ space where $k = F_0$.

- k can be very large

- k unknown

- data structure may not be fast enough.

But we can get exact answer.

Can show that one cannot get $(1-\varepsilon)$ approximation deterministically with sub-linear space.

But with randomization can get as a $(1-\varepsilon)$ approx estimate with probability $(1-\delta)$ in

$$O\left(\frac{1}{\varepsilon^2} \log \frac{1}{\delta}\right) \cdot polylog(n) \text{ space.}$$

## How?

- An old classical result from 1985 based on hashing and refined over years.

- A very recent surprisingly simple one based on sampling.

# Hashing based Algorithm

We will assume we have access
to an "ideal" hash function

$$h : [n] \longrightarrow [0,1]$$
$$\nearrow \text{ real interval.}$$

You can assume that instead of
real interval we are hashing to
$$[n^3] \text{ or some large } n^c$$

## Distinct Elements Hashing

- Pick random hash function $h : [n] \rightarrow [0,1]$

- $z \leftarrow \infty$

- While (stream is not done)
$$z \leftarrow \min(h(e_i), z).$$

- Output $\frac{1}{z} - 1$.

Memory: Only one number!

Lemma: Suppose $X_1, X_2, \ldots, X_k$ are independent random variables that are uniformly distributed in $[0,1]$. Let $Y = \min_i X_i$.

Then (i) $E[Y] = \dfrac{1}{k+1}$

(ii) $\text{Var}(Y) = \dfrac{1}{(k+1)(k+2)}$

Proof: Let $F_Y$ and $f_Y$ be the cdf and pdf of $Y$.

Then it is easy to see that

$$F_Y(t) = 0 \quad t \leq 0 \qquad F_Y(t) = 1 \quad t \geq 1$$

$$F_Y(t) = 1 - (1-t)^k \quad \text{for } t \text{ in } (0,1).$$

$$f_y(t) = k(1-t)^{k-1}$$

$$E[Y] = \int_0^1 (1-t)^k \, dt = -\frac{(1-t)^{k+1}}{k+1} \Big|_0^1$$

$$= \frac{1}{k+1} = \int_0^1 t \, k \, (1-t)^{k-1} \, dt$$

$$E[Y^2] = \int_0^1 t^2 \, k \, (1-t)^{k-1} \, dt$$

$$= \frac{2}{(k+1)(k+2)}$$

$$Var(Y) = \frac{2}{(k+1)(k+2)} - \frac{1}{(k+1)^2}$$

$$= \frac{2k}{(k+1)^2(k+2)} \approx \frac{2}{(k+1)(k+2)}$$

$\square$.

Thus the algorithm has an exact estimator for $F_0$ but variance is too high to get a good approximation. What does this mean?

Lemma: Suppose $F_0 = k$.

For $\frac{1}{Z} - 1$ to be with $(1 \pm \varepsilon)$ of $k$ we need $Z$ to be within $\frac{O(\varepsilon)}{k+1}$ of true value $\frac{1}{k+1}$.

But if we use Markov or Chebyshev one sees that it does not work.

# Variance reduction

When we have an estimator exact $Y$ for a quantity with true value $\alpha$ and $\mathrm{Var}(Y)$ is large we can obtain another estimator $Y'$ s.t $E[Y'] = \alpha$ and $\mathrm{Var}(Y')$ is $\leq \frac{1}{n} \mathrm{Var}(Y)$ by averaging independent estimators.

claim: Let $Y_1, Y_2, \ldots, Y_n$ be independent and let $Y = \frac{1}{n} \sum_{i=1}^{n} Y_i$

$$E[Y] = \frac{1}{n} \sum_{i} E[Y_i]$$

and $\mathrm{Var}(Y) = \frac{1}{n} \sum_{i=1}^{n} \mathrm{Var}(Y_i).$

Hence if each $Y_i$'s are iid with mean $\mu$ and variance $\sigma^2$

Then $E[Y] = \mu$ and $Var(Y) = \frac{\sigma^2}{n}$.

How can we apply this to $F_0$ estimation?

We have an exact estimator

$\mu = k.$ $\qquad Var(Z) = \frac{2}{(k+1)(k+2)}$.

Want an estimator $Z'$ such that $P_r\left[(1-\varepsilon)k \leq Z' \leq (1+\varepsilon)k\right]$

$\geq \frac{9}{10}$. Say.

$\varepsilon \in (0, \frac{1}{2})$ given.

- Run basic algorithm $h$ times independently and in parallel

- Let $Z_1, Z_2, \ldots, Z_h$ be the min values.

- Let $Z = \dfrac{1}{h} \sum\limits_{j=1}^{h} Z_j$

- Output $Y = \dfrac{1}{Z} - 1$

Claim: $E[Z] = \dfrac{1}{K+1}$

Claim: $\mathrm{Var}(Z) \leq \dfrac{2}{h(K+1)(K+2)}$.

Via Chebyshev's inequality

$$P_2\left[\left|Z - \frac{1}{k+1}\right| \geq \frac{\varepsilon}{3(k+1)}\right] \leq \frac{\frac{2}{4(k+1)(k+1)}}{\frac{\varepsilon^2}{9(k+1)^2}}$$

$$\leq \frac{18}{h\varepsilon^2}.$$

Thus if $\quad \dfrac{18}{h\varepsilon^2} \leq \dfrac{1}{10}$

$$\Rightarrow \quad Z \in \left(1 \pm \frac{\varepsilon}{3}\right)\frac{1}{k+1}$$

$$\Rightarrow \quad \frac{1}{Z} - 1 \in (1 \pm \varepsilon)\, k.$$

$$\Rightarrow \quad h \geq \frac{180}{\varepsilon^2}.$$

**Exercise:** To get $(1 \pm \varepsilon)$ approx with probability $\geq (1 - \delta)$ via Chebyshev argue that

$$h = \Omega\left(\frac{1}{\varepsilon^2} \cdot \frac{1}{\delta}\right) \text{ suffices.}$$

Space usage is $h$ times the space usage for single estimator.

Can we do better?

## Median Trick

Can use only $h = O\left(\frac{1}{\varepsilon^2} \ln \frac{1}{\delta}\right)$ estimators. How?

Suppose $Y$ is an estimator
for $k$ such that

$$P_r\left[(1-\varepsilon) \leq Y \leq (1+\varepsilon)k\right] \geq \frac{3}{4}$$

We can obtain such an
estimator using $O\left(\frac{1}{\varepsilon^2}\right)$ independent
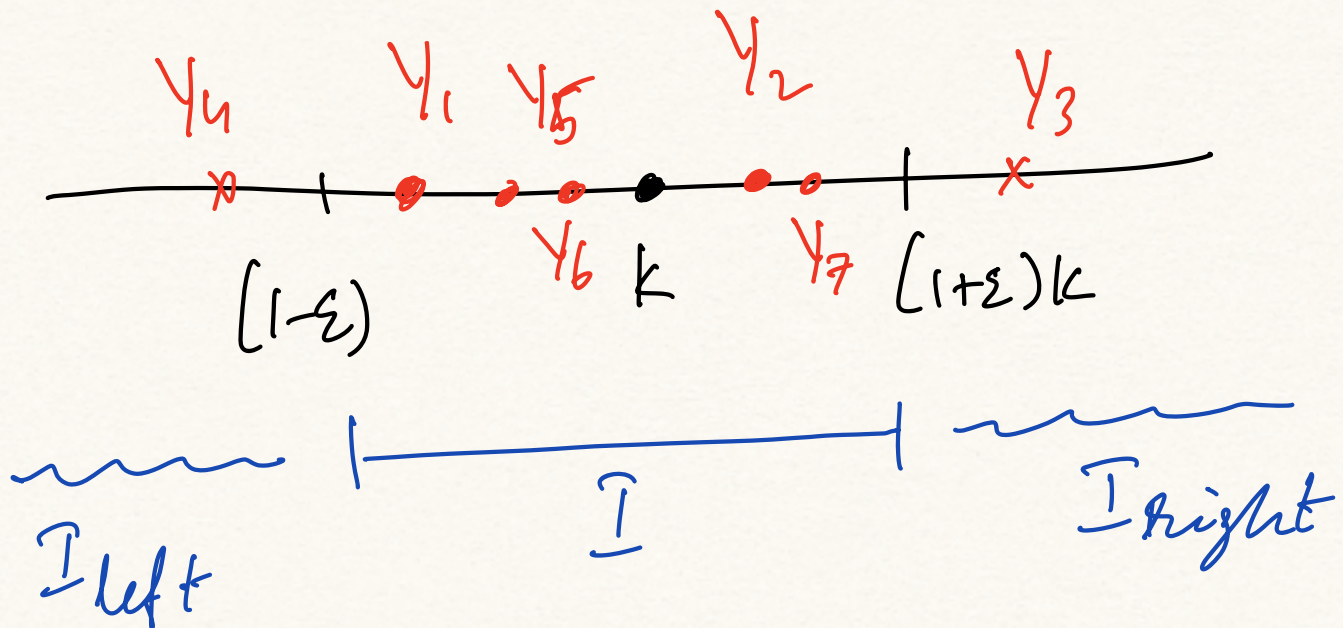values and averaging to reduce
variance and using Chebyshev as
above.

Let $Y_1, Y_2, ..., Y_\ell$ be independent
such estimators.

Lemma: Let $Y = \text{median}(Y_1, .., Y_\ell)$

Then $P_r\left[(1-\varepsilon)k \leq Y \leq (1+\varepsilon)k\right] \geq 1-\delta$

if $l = \Omega\left(\log\frac{1}{\delta}\right)$.

**Proof:** Use Chernoff bounds.



$A_i$ is event that $Y_i \in I$.

$$\Pr[A_i] \geq \frac{3}{4}.$$

For $Y$ to be outside $I$ it must be that $\geq \frac{l}{2}$ fall to the right of $I$ or $> \frac{l}{2}$ fall to the left of $I$.

Let $R_i$ be indicator for $Y_i \in I_{right}$

$$\Pr[R_i] \leq \frac{1}{4} \cdot$$

$$R = \sum_{i=1}^{l} R_i \qquad E[R] \leq \frac{l}{4} \cdot$$

$$\Pr\left[R \geq \frac{l}{2}\right] \leq e^{-c\frac{l}{4}}$$

$$\leq \frac{\delta}{2} \quad \text{if} \quad l = \Omega\left(\ln\frac{1}{\delta}\right)$$

Thus

$$\Pr\left[median\left(Y_1,.., Y_l\right) \in I_{right}\right] \leq \frac{\delta}{2}$$

By symmetric argument

$$\Pr\left[median\left(Y_1,.., Y_l\right) \in I_{left}\right] \leq \frac{\delta}{2}$$

$$\Rightarrow \Pr\left[median\left(Y_1,.., Y_l\right) \notin I\right] \leq \delta \cdot$$

$\square$

# Final alg:

Given $\varepsilon, \delta$.

- for $i = 1$ to $\ell$ do

  $X_i \leftarrow 0$

  for $j = 1$ to $h$ do

  Run DEHashing to obtain $Z_j$

  $X_i = Z_i + X_i$

  end for

  $X_i \leftarrow \frac{1}{h} X_i$

  $Y_i = \frac{1}{X_i} - 1$

- end for
- Output median $(Y_1, \cdots, Y_\ell)$.

<u>Space usage:</u> $\Omega\left(\frac{1}{\varepsilon^2}\log\frac{1}{\delta}\right)$ times space usage of basic estimator which is one number.

<u>Guarantee:</u> Output is a $(1\pm\varepsilon)$ approx with probability $\geqslant (1-\delta)$.

We abstract out what we did above.

<u>Defn:</u> An $(\varepsilon,\delta)$ randomized scheme for a non-negative quantity $\alpha$ is a randomized algorithm that outputs a value $X$ such that
$$\Pr\left[(1-\varepsilon)\alpha \leq X \leq (1+\varepsilon)X\right]\geqslant 1-\delta.$$

**Lemma:** Suppose we have a randomized algorithm that outputs $X \geq 0$ $E[X] = \alpha$ and $Var(X) \leq \beta \cdot \alpha^2$. Then one can use the algorithm $O\left(\frac{\beta}{\varepsilon^2} \ln \frac{1}{\delta}\right)$ time independently to obtain a $(\varepsilon, \delta)$-scheme.

From ideal hashing to "actual" hashing.

In streaming we are optimizing space so also need to look at space requirements of hash function.

- Can use pairwise independent hash function.

- no real interval but can approximate $[0,1]$ by $[0, \frac{1}{n^3}, \frac{2}{n^3}, \ldots, 1]$ hence $h: [n] \to [n^3]$ would be ok.

- still need small tweak to the algorithm since we are using limited independence.

New simple algorithm due to [CMV22]

Idea:

Suppose we sample each element of the stream with some fixed probability $p$

Let $B$ be the sample.

If $m$ is the length of the stream then

$$E[|B|] = pm$$

$\Rightarrow \dfrac{|B|}{p}$ is an unbiased estimator of $m$.

moreover by Chernoff bounds if $E\left[\dfrac{|B|}{p}\right] \geqslant \dfrac{c}{\varepsilon^2} \ln \dfrac{1}{\delta}$ we can

Show that $\frac{|B|}{p}$ is a $(1 \pm \varepsilon)$ approx for $m$ with prob $\geq 1 - \delta$.

But we don't want to estimate $m$ but want to estimate $F_0$.

- Is there a way to sample only distinct elements?

- How do we choose $p$ so that $\frac{|B|}{p}$ is sufficiently large but not too large for otherwise we will need to store $|B|$.

We will address first issue.
Suppose we maintain a sample
$B$ of the stream

$$B \leftarrow \emptyset \quad, \quad m \leftarrow 0$$
While (stream is not done)

$$m \leftarrow m+1$$
$c_m$ current element
add $c_m$ to $B$ with prob $p$.

How do we make $B$ sample only
distinct elements?

Trick. If $c_m \in B$ remove it
since we have detected a duplicate
and then sample the new element
with prob $p$.

distinct-element-sample (p)

$B \leftarrow \phi$ , $m \leftarrow 0$
while (stream is not done)

$m \leftarrow m+1$
$e_m$ current element
If $e_m \in B$, $B \leftarrow B - \{e_m\}$
Add $e_m$ to $B$ with prob $p$.

Output $\dfrac{|B|}{p}$ as estimate for $F_0$.

Second issue: how do we set $p$?

If we set $p$ too low, $|B|$ will be too small to get accurate estimate. If we set $p$ too high

$|B|$ will be too large.

Key idea is to "guess" $F_0$ and reduce $p$ if $|B|$ gets too large.

distin-elements-Sample $(n, \varepsilon)$

$\quad p \leftarrow 1, \quad B \leftarrow \phi, \quad m \leftarrow \phi$

$\quad \tau \leftarrow c \dfrac{\log n}{\varepsilon^2}$ for sufficiently large $C$.

$\quad$ While (stream is not empty) do

$\qquad m \leftarrow m+1$

$\qquad e_m$ current item

$\qquad B \leftarrow B - \{e_m\}$

$\qquad$ With prob $p$:

$\qquad \quad$ add $e_m$ to $B$

$\qquad$ if $|B| \geq \tau$

$\qquad \quad$ discard each $e_i \in B$

$\qquad \qquad$ with prob $\frac{1}{2}$.

end while

Output $\dfrac{|B|}{p}$

Theorem: The algorithm used $O\left(\dfrac{\log n}{\varepsilon^2}\right)$ words of memory and outputs an estimate that is within $(1\pm\varepsilon)F_0$ with high probability.