Topics in TCS

Counting distinct elements

Raphaël Clifford







10 distinct colours

Naive counting solution

Sort the elements. $O(m \log m)$ time.

Naive counting solution

Sort the elements. $O(m \log m)$ time.

Traverse from left to right. O(m) time.

Naive counting solution

Sort the elements. $O(m \log m)$ time.

Traverse from left to right. O(m) time.

 $O(m \log m)$ time overall BUT O(m) words of space.

Naive counting solution

Sort the elements. $O(m \log m)$ time.

Traverse from left to right. O(m) time.

 $O(m \log m)$ time overall BUT O(m) words of space.

NOT one-pass.



Slightly less naive counting solution

Use a balanced binary search tree.







At the conclusion the number of distinct values will be the number of nodes in tree.



At the conclusion the number of distinct values will be the number of nodes in tree.

Running time $O(\log m)$ per FIND and INSERT operation making $O(m \log m)$ time overall BUT O(m) words of space.



At the conclusion the number of distinct values will be the number of nodes in tree.

Running time $O(\log m)$ per FIND and INSERT operation making $O(m \log m)$ time overall BUT O(m) words of space.

One-pass \checkmark



At the conclusion the number of distinct values will be the number of nodes in tree.

Running time $O(\log m)$ per FIND and INSERT operation making $O(m \log m)$ time overall BUT O(m) words of space.

One-pass \checkmark

No deterministic sub-linear space one-pass solution is possible.



At the conclusion the number of distinct values will be the number of nodes in tree.

Running time $O(\log m)$ per FIND and INSERT operation making $O(m \log m)$ time overall BUT O(m) words of space.

One-pass √

No deterministic sub-linear space one-pass solution is possible.

We will need a randomised algorithm.

```
stream \langle a_1, \ldots, a_m \rangle, a_i \in [n]
initialise
choose random h: [n] \rightarrow [n]
SIMPLE-COUNT
Set M = h(a_1)
For each i \geq 2
   if h(a_i) < M
          set M = h(a_i)
return \hat{d} = n/M
```

```
stream \langle a_1, \ldots, a_m \rangle, a_i \in [n]
initialise
choose random h: [n] \rightarrow [n]
SIMPLE-COUNT
Set M = h(a_1)
For each i \ge 2
   if h(a_i) < M
          set M = h(a_i)
return \hat{d} = n/M
```

```
stream \langle a_1, \ldots, a_m 
angle, a_i \in [n]
initialise
choose random h: [n] \rightarrow [n]
SIMPLE-COUNT
Set M = h(a_1)
For each i \ge 2
   if h(a_i) < M
          set M = h(a_i)
return \hat{d} = n/M
```

Worked example

Stream ⟨5, 4, 5, 7, 4, 5, 7, 4⟩.
 n = 10.

```
stream \langle a_1, \ldots, a_m \rangle, a_i \in [n]
initialise
choose random h: [n] \rightarrow [n]
SIMPLE-COUNT
Set M = h(a_1)
For each i \ge 2
   if h(a_i) < M
          set M = h(a_i)
return \hat{d} = n/M
```

- Stream $\langle 5, 4, 5, 7, 4, 5, 7, 4 \rangle$. n = 10.
- Randomly hash to $\langle 10, 9, 10, 6, 9, 10, 6, 9 \rangle$

```
stream \langle a_1, \ldots, a_m \rangle, a_i \in [n]
initialise
choose random h: [n] \rightarrow [n]
SIMPLE-COUNT
Set M = h(a_1)
For each i \ge 2
   if h(a_i) < M
          set M = h(a_i)
return \hat{d} = n/M
```

- Stream $\langle 5, 4, 5, 7, 4, 5, 7, 4 \rangle$. n = 10.
- Randomly hash to $\langle 10, 9, 10, 6, 9, 10, 6, 9 \rangle$
- M = 6 at termination.

```
stream \langle a_1,\ldots,a_m
angle,a_i\in[n]
initialise
choose random h: [n] \rightarrow [n]
SIMPLE-COUNT
Set M = h(a_1)
For each i \ge 2
   if h(a_i) < M
         set M = h(a_i)
return \hat{d} = n/M
```

- Stream $\langle 5, 4, 5, 7, 4, 5, 7, 4 \rangle$. n = 10.
- Randomly hash to $\langle 10, 9, 10, 6, 9, 10, 6, 9 \rangle$
- *M* = 6 at termination.
- Return $\hat{d} = 10/6 = 1\frac{2}{3}$

```
stream \langle a_1, \dots, a_m 
angle , a_i \in [n]
initialise
choose random h: [n] \rightarrow [n]
SIMPLE-COUNT
Set M = h(a_1)
For each i \ge 2
   if h(a_i) < M
          set M = h(a_i)
return \hat{d} = n/M
```

- Stream $\langle 5, 4, 5, 7, 4, 5, 7, 4 \rangle$. n = 10.
- Randomly hash to $\langle 10, 9, 10, 6, 9, 10, 6, 9 \rangle$
- *M* = 6 at termination.
- Return $\hat{d} = 10/6 = 1\frac{2}{3}$
- True answer is 3.

```
stream \langle a_1,\ldots,a_m
angle,a_i\in[n]
initialise
choose random h: [n] \rightarrow [n]
SIMPLE-COUNT
Set M = h(a_1)
For each i \ge 2
   if h(a_i) < M
         set M = h(a_i)
return \hat{d} = n/M
```

Worked example. 2nd try

- Stream (5, 4, 5, 7, 4, 5, 7, 4). n = 10.
- Rehash: $\langle 4, 3, 4, 8, 3, 4, 8, 3 \rangle$

• True answer is 3.

```
stream \langle a_1, \ldots, a_m 
angle, a_i \in [n]
initialise
choose random h: [n] \rightarrow [n]
SIMPLE-COUNT
Set M = h(a_1)
For each i \ge 2
   if h(a_i) < M
          set M = h(a_i)
return \hat{d} = n/M
```

Worked example. 2nd try

- Stream $\langle 5, 4, 5, 7, 4, 5, 7, 4 \rangle$. n = 10.
- Rehash: $\langle 4, 3, 4, 8, 3, 4, 8, 3 \rangle$

• True answer is 3.

```
stream \langle a_1, \dots, a_m 
angle , a_i \in [n]
initialise
choose random h: [n] \rightarrow [n]
SIMPLE-COUNT
Set M = h(a_1)
For each i \ge 2
   if h(a_i) < M
          set M = h(a_i)
return \hat{d} = n/M
```

Worked example. 2nd try

- Stream $\langle 5, 4, 5, 7, 4, 5, 7, 4 \rangle$. n = 10.
- Rehash: $\langle 4, 3, 4, 8, 3, 4, 8, 3 \rangle$
- M = 3 at termination.
- Return $\hat{d} = 10/3 = 3\frac{1}{3}$
- True answer is 3.

Lemma

Let $(Z_1, ..., Z_N)$ be an array of pairwise independent indicator random variables with $Pr(Z_i = 1) = p$ and let $W = \sum_{i=1}^N Z_i$, then $\mathbb{E}(W) = Np$, var(W) = Np(1-p) and

$$\Pr(W > 0) \le Np$$
 and $\Pr(W = 0) \le \frac{1}{Np}$ (1)

Proof.

Lemma

Let $(Z_1, ..., Z_N)$ be an array of pairwise independent indicator random variables with $Pr(Z_i = 1) = p$ and let $W = \sum_{i=1}^N Z_i$, then $\mathbb{E}(W) = Np$, var(W) = Np(1-p) and

$$\Pr(W > 0) \le Np$$
 and $\Pr(W = 0) \le \frac{1}{Np}$ (1)

Proof.

► E(W) = Np, var(W) = Np(1 - p) both follow from linearity of expectation and the preliminary probability notes.

Lemma

Let $(Z_1, ..., Z_N)$ be an array of pairwise independent indicator random variables with $Pr(Z_i = 1) = p$ and let $W = \sum_{i=1}^N Z_i$, then $\mathbb{E}(W) = Np$, var(W) = Np(1-p) and

$$\Pr(W > 0) \le Np$$
 and $\Pr(W = 0) \le \frac{1}{Np}$ (1)

Proof.

- ► 𝔅(𝒜) = №𝑘, var(𝒜) = №𝑘(1 − 𝑘) both follow from linearity of expectation and the preliminary probability notes.
- $Pr(W > 0) \le Np$ is by the union bound and that $Pr(Z_i = 1) = p$.

Lemma

Let $(Z_1, ..., Z_N)$ be an array of pairwise independent indicator random variables with $Pr(Z_i = 1) = p$ and let $W = \sum_{i=1}^N Z_i$, then $\mathbb{E}(W) = Np$, var(W) = Np(1-p) and

$$\Pr(W > 0) \le Np$$
 and $\Pr(W = 0) \le \frac{1}{Np}$ (1)

Proof.

- ► 𝔅(𝒴) = Np, var(𝒴) = Np(1 − p) both follow from linearity of expectation and the preliminary probability notes.
- $Pr(W > 0) \le Np$ is by the union bound and that $Pr(Z_i = 1) = p$.

►
$$\Pr(W = 0) \leq \Pr(|W - \mathbb{E}(W)| \geq \mathbb{E}(W)) \leq \frac{\operatorname{var}(W)}{(Np)^2} \leq \frac{1}{Np}$$

$\label{eq:SIMPLE-Count analysis} SIMPLE-COUNT analysis$

• Let us consider the set of distinct values in the stream S with d = |S|.

- Let us consider the set of distinct values in the stream S with d = |S|.
- Let r.v. Y_a be the number of elements that are hashed to a value less than or equal to a, for arbitrary a.

- Let us consider the set of distinct values in the stream S with d = |S|.
- Let r.v. Y_a be the number of elements that are hashed to a value less than or equal to a, for arbitrary a.
- Y_a corresponds to W from (1) with N = d and p = a/n. Therefore:

$$\Pr(Y_a > 0) = \Pr(M \le a) \le da/n$$
 and $\Pr(M > b) \le \frac{n}{db}$ (2)

- Let us consider the set of distinct values in the stream S with d = |S|.
- Let r.v. Y_a be the number of elements that are hashed to a value less than or equal to a, for arbitrary a.
- Y_a corresponds to W from (1) with N = d and p = a/n. Therefore:

$$\frac{\Pr(Y_a > 0) = \Pr(M \le a)}{(db)} \le da/n \quad \text{and} \quad \Pr(M > b) \le \frac{n}{db} \quad (2)$$
s an element less than

n

or equal to a iff the min is at most a

There i

- Let us consider the set of distinct values in the stream S with d = |S|.
- Let r.v. Y_a be the number of elements that are hashed to a value less than or equal to a, for arbitrary a.
- Y_a corresponds to W from (1) with N = d and p = a/n. Therefore:



- Let us consider the set of distinct values in the stream S with d = |S|.
- Let r.v. Y_a be the number of elements that are hashed to a value less than or equal to a, for arbitrary a.
- Y_a corresponds to W from (1) with N = d and p = a/n. Therefore:

$$\Pr(Y_a > 0) = \Pr(M \le a) \le da/n$$
 and $\Pr(M > b) \le \frac{n}{db}$ (2)

Now let da/n = 1/3 and n/db = 1/3, so that a = n/(3d) and b = 3n/d, then (2) becomes

$$\Pr(M \le n/(3d)) \le 1/3$$
 and $\Pr(M > 3n/d) \le 1/3$

or

$$\Pr(d \leq \hat{d}/3) \leq 1/3$$
 and $\Pr(d > 3\hat{d}) \leq 1/3$
using $\hat{d} = n/M.$
SIMPLE-COUNT space/time

$$\Pr(d \leq \hat{d}/3) \leq 1/3$$
 and $\Pr(d > 3\hat{d}) \leq 1/3$

```
\begin{array}{l} \texttt{stream} \ \left\langle \textbf{\textit{a}}_1, \ldots, \textbf{\textit{a}}_m \right\rangle, \textbf{\textit{a}}_i \in [n] \\ \texttt{initialise} \end{array}
Choose random h: [n] \rightarrow [n]
SIMPLE-COUNT
Set M = h(a_1)
For each i \geq 2
     if h(a_i) < M
                  set M = h(a_i)
return \hat{d} = n/M
```

$\operatorname{SIMPLE-COUNT}$ space/time

Our random estimate \hat{d} gives us:

$$\Pr(d \leq \hat{d}/3) \leq 1/3$$
 and $\Pr(d > 3\hat{d}) \leq 1/3$

```
\texttt{stream} \hspace{0.1 in} \langle \textbf{\textit{a}}_1, \dots, \textbf{\textit{a}}_m \rangle \textbf{,} \textbf{\textit{a}}_i \in [\textbf{\textit{n}}] \\ \texttt{initialise}
Choose random h:[n] \to [n]
SIMPLE-COUNT
Set M = h(a_1)
For each i \ge 2
     if h(a_i) < M
                   set M = h(a_i)
return \hat{d} = n/M
```

• Running time O(m)

$\operatorname{SIMPLE-COUNT}$ space/time

$$\Pr(d \leq \hat{d}/3) \leq 1/3$$
 and $\Pr(d > 3\hat{d}) \leq 1/3$

```
\texttt{stream} \hspace{0.1 in} \langle \textbf{\textit{a}}_1, \dots, \textbf{\textit{a}}_m \rangle \textbf{,} \textbf{\textit{a}}_i \in [\textbf{\textit{n}}] \\ \texttt{initialise}
Choose random h: [n] \rightarrow [n]
SIMPLE-COUNT
Set M = h(a_1)
For each i \ge 2
      if h(a_i) < M
                    set M = h(a_i)
return \hat{d} = n/M
```

- Running time O(m)
- One-pass √

${\rm SIMPLE}\text{-}{\rm COUNT} \text{ space}/\text{time}$

$$\Pr(d \leq \hat{d}/3) \leq 1/3$$
 and $\Pr(d > 3\hat{d}) \leq 1/3$

```
\begin{array}{l} \texttt{stream} \ \left\langle \textbf{\textit{a}}_1, \ldots, \textbf{\textit{a}}_m \right\rangle, \textbf{\textit{a}}_i \in [\textbf{\textit{n}}] \\ \texttt{initialise} \end{array}
Choose random h: [n] \rightarrow [n]
SIMPLE-COUNT
Set M = h(a_1)
For each i \ge 2
       if h(a_i) < M
                     set M = h(a_i)
return \hat{d} = n/M
```

- Running time O(m)
- One-pass √
- Space O(log n) √

${\rm SIMPLE}\text{-}{\rm COUNT} \text{ space}/\text{time}$

$$\Pr(d \leq \hat{d}/3) \leq 1/3$$
 and $\Pr(d > 3\hat{d}) \leq 1/3$

```
\begin{array}{l} \texttt{stream} \ \left\langle \textbf{\textit{a}}_1, \ldots, \textbf{\textit{a}}_m \right\rangle, \textbf{\textit{a}}_i \in [\textbf{\textit{n}}] \\ \texttt{initialise} \end{array}
Choose random h: [n] \rightarrow [n]
SIMPLE-COUNT
Set M = h(a_1)
For each i \ge 2
      if h(a_i) < M
                     set M = h(a_i)
return \hat{d} = n/M
```

- Running time O(m)
- One-pass √
- Space O(log n) √
- Can we use less space? (Sort of)

${\rm SIMPLE}\text{-}{\rm COUNT} \text{ space}/\text{time}$

$$\Pr(d \leq \hat{d}/3) \leq 1/3$$
 and $\Pr(d > 3\hat{d}) \leq 1/3$

```
stream \langle a_1, \ldots, a_m \rangle, a_i \in [n] initialise
Choose random h: [n] \rightarrow [n]
SIMPLE-COUNT
Set M = h(a_1)
For each i \ge 2
    if h(a_i) < M
            set M = h(a_i)
return \hat{d} = n/M
```

- Running time O(m)
- One-pass √
- Space O(log n) √
- Can we use less space? (Sort of)
- The error probability is pretty bad. Can we fix that? (Yes)

```
initialise
choose random h: [n] \rightarrow [n]
set z = 0
TIDEMARK(a_i)
if ZEROS(h(a_i)) > z
set z = ZEROS(h(a_i))
OUTPUT 2^{z+\frac{1}{2}}
```

```
initialise

choose random h: [n] \rightarrow [n]

set z = 0

TIDEMARK(a_i)

if ZEROS(h(a_i)) > z

set z = ZEROS(h(a_i))

OUTPUT 2^{z+\frac{1}{2}}
```

h chosen from a pairwise random family of hash functions.

```
initialise

choose random h: [n] \rightarrow [n]

set z = 0

TIDEMARK(a_i)

if ZEROS(h(a_i)) > z

set z = ZEROS(h(a_i))

OUTPUT 2^{z+\frac{1}{2}}
```

- *h* chosen from a pairwise random family of hash functions.
- ZEROS counts the number of trailing zeros in the binary representation of a positive integer.

```
initialise

choose random h: [n] \rightarrow [n]

set z = 0

TIDEMARK(a_i)

if ZEROS(h(a_i)) > z

set z = ZEROS(h(a_i))

OUTPUT 2^{z+\frac{1}{2}}
```

- *h* chosen from a pairwise random family of hash functions.
- ZEROS counts the number of trailing zeros in the binary representation of a positive integer.
- Finds the maximum value of ZEROS(*h*(*a_i*)) overall.

```
initialise

choose random h: [n] \rightarrow [n]

set z = 0

TIDEMARK(a_i)

if ZEROS(h(a_i)) > z

set z = ZEROS(h(a_i))

OUTPUT 2^{z+\frac{1}{2}}
```

- h chosen from a pairwise random family of hash functions.
- ZEROS counts the number of trailing zeros in the binary representation of a positive integer.
- Finds the maximum value of ZEROS(*h*(*a_i*)) overall.

Let's try it. Stream (5, 4, 5, 7, 4, 5, 7, 4). n = 10.

```
initialise

choose random h: [n] \rightarrow [n]

set z = 0

TIDEMARK(a_i)

if ZEROS(h(a_i)) > z

set z = ZEROS(h(a_i))

OUTPUT 2^{z+\frac{1}{2}}
```

- *h* chosen from a pairwise random family of hash functions.
- ZEROS counts the number of trailing zeros in the binary representation of a positive integer.
- Finds the maximum value of ZEROS(*h*(*a_i*)) overall.

Let's try it. Stream (5, 4, 5, 7, 4, 5, 7, 4). n = 10.

• Hashed to $\langle 8, 2, 8, 1, 2, 8, 1, 2 \rangle$

```
initialise

choose random h: [n] \rightarrow [n]

set z = 0

TIDEMARK(a_i)

if ZEROS(h(a_i)) > z

set z = ZEROS(h(a_i))

OUTPUT 2^{z+\frac{1}{2}}
```

- *h* chosen from a pairwise random family of hash functions.
- ZEROS counts the number of trailing zeros in the binary representation of a positive integer.
- Finds the maximum value of ZEROS(*h*(*a_i*)) overall.

Let's try it. Stream (5, 4, 5, 7, 4, 5, 7, 4). n = 10.

- Hashed to $\langle 8, 2, 8, 1, 2, 8, 1, 2 \rangle$
- ▶ In binary: $\langle 1000, 10, 1000, 1, 10, 1000, 1, 10 \rangle$

```
initialise

choose random h: [n] \rightarrow [n]

set z = 0

TIDEMARK(a_i)

if ZEROS(h(a_i)) > z

set z = ZEROS(h(a_i))

OUTPUT 2^{z+\frac{1}{2}}
```

- *h* chosen from a pairwise random family of hash functions.
- ZEROS counts the number of trailing zeros in the binary representation of a positive integer.
- Finds the maximum value of ZEROS(*h*(*a_i*)) overall.

Let's try it. Stream (5, 4, 5, 7, 4, 5, 7, 4). n = 10.

- Hashed to $\langle 8, 2, 8, 1, 2, 8, 1, 2 \rangle$
- ▶ In binary: (1000, 10, 1000, 1, 10, 1000, 1, 10)
- Max value of ZEROS is 3. Return $2^{3+1/2} \approx 11.3$

```
initialise

choose random h: [n] \rightarrow [n]

set z = 0

TIDEMARK(a_i)

if ZEROS(h(a_i)) > z

set z = ZEROS(h(a_i))

OUTPUT 2^{z+\frac{1}{2}}
```

- *h* chosen from a pairwise random family of hash functions.
- ZEROS counts the number of trailing zeros in the binary representation of a positive integer.
- Finds the maximum value of ZEROS(*h*(*a_i*)) overall.

Let's try it. Stream (5, 4, 5, 7, 4, 5, 7, 4). n = 10.

- Hashed to $\langle 8, 2, 8, 1, 2, 8, 1, 2 \rangle$
- ▶ In binary: (1000, 10, 1000, 1, 10, 1000, 1, 10)
- Max value of $Z_{\rm EROS}$ is 3. Return $2^{3+1/2} \approx 11.3$

True value 3. What happens if we pick another hash function?

```
initialise

choose random h: [n] \rightarrow [n]

set z = 0

TIDEMARK(a_i)

if ZEROS(h(a_i)) > z

set z = ZEROS(h(a_i))

OUTPUT 2^{z+\frac{1}{2}}
```

- *h* chosen from a pairwise random family of hash functions.
- ZEROS counts the number of trailing zeros in the binary representation of a positive integer.
- Finds the maximum value of ZEROS(*h*(*a_i*)) overall.

Let's try it. Stream (5, 4, 5, 7, 4, 5, 7, 4). n = 10. Hashed to (4, 3, 4, 1, 3, 4, 1, 3)

```
initialise

choose random h: [n] \rightarrow [n]

set z = 0

TIDEMARK(a_i)

if ZEROS(h(a_i)) > z

set z = ZEROS(h(a_i))

OUTPUT 2^{z+\frac{1}{2}}
```

- *h* chosen from a pairwise random family of hash functions.
- ZEROS counts the number of trailing zeros in the binary representation of a positive integer.
- Finds the maximum value of ZEROS(*h*(*a_i*)) overall.

Let's try it. Stream (5, 4, 5, 7, 4, 5, 7, 4). n = 10.

- Hashed to $\langle 4, 3, 4, 1, 3, 4, 1, 3 \rangle$
- In binary: $\langle 100, 11, 100, 1, 11, 100, 1, 11 \rangle$

```
initialise

choose random h: [n] \rightarrow [n]

set z = 0

TIDEMARK(a_i)

if ZEROS(h(a_i)) > z

set z = ZEROS(h(a_i))

OUTPUT 2^{z+\frac{1}{2}}
```

- *h* chosen from a pairwise random family of hash functions.
- ZEROS counts the number of trailing zeros in the binary representation of a positive integer.
- Finds the maximum value of ZEROS(*h*(*a_i*)) overall.

Let's try it. Stream (5, 4, 5, 7, 4, 5, 7, 4). n = 10.

- Hashed to $\langle 4, 3, 4, 1, 3, 4, 1, 3 \rangle$
- \blacktriangleright In binary: $\langle 100, 11, 100, 1, 11, 100, 1, 11 \rangle$
- Max value of ZEROS is 2. Return $2^{2+1/2} \approx 5.7$. We were luckier.

Understanding the TIDEMARK algorithm

- ZEROS(x) = 0 $x = 1, 3, 5, \dots$ ZEROS $(x) \ge 1$ $x = 2, 4, 6, \ldots$ $Z_{EROS}(x) \ge 2$ x = 4, 8, 12, ... $Z_{EROS}(x) \geq 3$ $\operatorname{ZEROS}(x) \geq 4$ $x = 16, 32, 48, \ldots$
 - $x = 8, 16, 24, \ldots$

$\operatorname{Zeros}(x) = 0$	$x=1,3,5,\ldots$	p = 1/2
$\operatorname{Zeros}(x) \geq 1$	$x = 2, 4, 6, \dots$	p = 1/2
$\operatorname{Zeros}(x) \geq 2$	$x = 4, 8, 12, \dots$	p=1/4
$\operatorname{Zeros}(x) \geq 3$	$x = 8, 16, 24, \dots$	p=1/8
$\operatorname{Zeros}(x) \geq 4$	$x = 16, 32, 48, \dots$	p = 1/16

0 +

1

$\operatorname{ZEROS}(x) = 0$	$x = 1, 3, 5, \ldots$	p=1/2
$\operatorname{Zeros}(x) \geq 1$	$x = 2, 4, 6, \dots$	p = 1/2
$\operatorname{Zeros}(x) \geq 2$	$x = 4, 8, 12, \dots$	p=1/4
$\operatorname{Zeros}(x) \geq 3$	$x = 8, 16, 24, \dots$	p=1/8
$\operatorname{Zeros}(x) \geq 4$	$x = 16, 32, 48, \dots$	p = 1/16



$\operatorname{ZEROS}(x) = 0$	$x = 1, 3, 5, \ldots$	p=1/2
$\operatorname{Zeros}(x) \geq 1$	$x = 2, 4, 6, \dots$	p=1/2
$\operatorname{Zeros}(x) \geq 2$	$x = 4, 8, 12, \dots$	p=1/4
$\operatorname{Zeros}(x) \geq 3$	$x = 8, 16, 24, \dots$	p=1/8
$\operatorname{ZEROS}(x) \geq 4$	$x = 16, 32, 48, \dots$	p=1/16



$\operatorname{ZEROS}(x) = 0$	$x = 1, 3, 5, \ldots$	p=1/2	
$\operatorname{Zeros}(x) \geq 1$	$x = 2, 4, 6, \dots$	p=1/2	
$\operatorname{Zeros}(x) \geq 2$	$x = 4, 8, 12, \dots$	p=1/4	
$\operatorname{Zeros}(x) \geq 3$	$x = 8, 16, 24, \dots$	p=1/8	
$\operatorname{ZEROS}(x) \geq 4$	$x = 16, 32, 48, \dots$	p = 1/16	



$\operatorname{ZEROS}(x) = 0$	$x=1,3,5,\ldots$	p=1/2	
$\operatorname{ZEROS}(x) \geq 1$	$x = 2, 4, 6, \ldots$	p=1/2	
$\operatorname{ZEROS}(x) \geq 2$	$x = 4, 8, 12, \dots$	p=1/4	
$\operatorname{Zeros}(x) \geq 3$	<i>x</i> = 8, 16, 24,	p=1/8	
$\operatorname{Zeros}(x) \geq 4$	$x = 16, 32, 48, \ldots$	p = 1/16	



For $j \in [n]$, define indicator r.v. $X_{r,j} = 1$ if $\operatorname{ZEROS}(h(j)) \ge r$.

For $j \in [n]$, define indicator r.v. $X_{r,j} = 1$ if $\operatorname{ZEROS}(h(j)) \ge r$.

Define r.v. $Y_r = \sum_{j, f_j > 0} X_{r,j}$. That is the number of hashed values in the stream that have at least r trailing zeros.

For $j \in [n]$, define indicator r.v. $X_{r,j} = 1$ if $\operatorname{ZEROS}(h(j)) \ge r$.

Define r.v. $Y_r = \sum_{j, f_j > 0} X_{r,j}$. That is the number of hashed values in the stream that have at least *r* trailing zeros.

Let T be the final value of z,

For $j \in [n]$, define indicator r.v. $X_{r,j} = 1$ if $\operatorname{ZEROS}(h(j)) \ge r$.

Define r.v. $Y_r = \sum_{j, f_j > 0} X_{r,j}$. That is the number of hashed values in the stream that have at least *r* trailing zeros.

Let T be the final value of z,

 $Y_r > 0$ iff $T \ge r$

For $j \in [n]$, define indicator r.v. $X_{r,j} = 1$ if $\operatorname{ZEROS}(h(j)) \ge r$.

Define r.v. $Y_r = \sum_{j, f_j > 0} X_{r,j}$. That is the number of hashed values in the stream that have at least *r* trailing zeros.

Let T be the final value of z,

$$Y_r > 0 ext{ iff } T \ge r$$

 $Y_r = 0 ext{ iff } T \le r - 1$

For $j \in [n]$, define indicator r.v. $X_{r,j} = 1$ if $\operatorname{ZEROS}(h(j)) \ge r$.

Define r.v. $Y_r = \sum_{j, f_j > 0} X_{r,j}$. That is the number of hashed values in the stream that have at least *r* trailing zeros.

Let T be the final value of z,

$$egin{aligned} &Y_r > 0 \ ext{iff} \ T \geq r \ &Y_r = 0 \ ext{iff} \ T \leq r-1 \ &\mathbb{E}(X_{r,j}) = \Pr(ext{zeros}(h(j)) \geq r) = \Pr(2^r \ ext{divides} \ h(j)) \end{aligned}$$

For $j \in [n]$, define indicator r.v. $X_{r,j} = 1$ if $\operatorname{ZEROS}(h(j)) \ge r$.

Define r.v. $Y_r = \sum_{j, f_j > 0} X_{r,j}$. That is the number of hashed values in the stream that have at least r trailing zeros.

Let T be the final value of z,

$$egin{aligned} &Y_r > 0 ext{ iff } T \geq r \ &Y_r = 0 ext{ iff } T \leq r-1 \ &\mathbb{E}(X_{r,j}) = \Pr(\operatorname{ZEROS}(h(j)) \geq r) = \Pr(2^r ext{ divides } h(j)) \end{aligned}$$

Assuming $2^r \leq n$ then

$$1/2^r \le \Pr(2^r \text{ divides } h(j)) \le 1/2^{r-1}$$

For $j \in [n]$, define indicator r.v. $X_{r,j} = 1$ if $\operatorname{ZEROS}(h(j)) \ge r$.

Define r.v. $Y_r = \sum_{j, f_j > 0} X_{r,j}$. That is the number of hashed values in the stream that have at least *r* trailing zeros.

Let T be the final value of z,

$$egin{aligned} &Y_r > 0 ext{ iff } T \geq r \ &Y_r = 0 ext{ iff } T \leq r-1 \ &\mathbb{E}(X_{r,j}) = \Pr(\operatorname{ZEROS}(h(j)) \geq r) = \Pr(2^r ext{ divides } h(j)) \end{aligned}$$

Assuming $2^r \leq n$ then

$$1/2^r \le \Pr(2^r \text{ divides } h(j)) \le 1/2^{r-1}$$

If *n* is a power of 2, $\mathbb{E}(X_{r,j}) = 1/2^r$

For simplicity, assume n is a power of 2.

$$\mathbb{E}(Y_r) = \sum_{j:f_j>0} \mathbb{E}(X_{r,j}) = \frac{d}{2^r}$$

For simplicity, assume n is a power of 2.

$$\mathbb{E}(Y_r) = \sum_{j:f_j > 0} \mathbb{E}(X_{r,j}) = \frac{d}{2^r}$$
$$\operatorname{var}(Y_r) = \sum_{j:f_j > 0} \operatorname{var}(X_{r,j}) \le \sum_{j:f_j > 0} \mathbb{E}(X_{r,j}^2) = \sum_{j:f_j > 0} \mathbb{E}(X_{r,j}) = \frac{d}{2^r}$$

For simplicity, assume n is a power of 2.

$$\mathbb{E}(Y_r) = \sum_{j:f_j > 0} \mathbb{E}(X_{r,j}) = \frac{d}{2^r}$$
$$\operatorname{var}(Y_r) = \sum_{j:f_j > 0} \operatorname{var}(X_{r,j}) \le \sum_{j:f_j > 0} \mathbb{E}(X_{r,j}^2) = \sum_{j:f_j > 0} \mathbb{E}(X_{r,j}) = \frac{d}{2^r}$$

By Markov's inequality,

$$\mathsf{Pr}(Y_r > 0) = \mathsf{Pr}(Y_r \ge 1) \le \frac{\mathbb{E}(Y_r)}{1} = \frac{d}{2^r}$$

For simplicity, assume n is a power of 2.

$$\mathbb{E}(Y_r) = \sum_{j:f_j > 0} \mathbb{E}(X_{r,j}) = \frac{d}{2^r}$$
$$\operatorname{var}(Y_r) = \sum_{j:f_j > 0} \operatorname{var}(X_{r,j}) \le \sum_{j:f_j > 0} \mathbb{E}(X_{r,j}^2) = \sum_{j:f_j > 0} \mathbb{E}(X_{r,j}) = \frac{d}{2^r}$$

By Markov's inequality,

$$\Pr(Y_r > 0) = \Pr(Y_r \ge 1) \le \frac{\mathbb{E}(Y_r)}{1} = \frac{d}{2^r}$$

By Chebyshev's inequality,

$$\Pr(Y_r=0) \leq \Pr(|Y_r - \mathbb{E}(Y_r)| \geq d/2^r) \leq \frac{\operatorname{var}(Y_r)}{(d/2^r)^2} \leq \frac{2^r}{d}$$
For simplicity, assume n is a power of 2.

$$\mathbb{E}(Y_r) = \sum_{j:f_j > 0} \mathbb{E}(X_{r,j}) = \frac{d}{2^r}$$
$$\operatorname{var}(Y_r) = \sum_{j:f_j > 0} \operatorname{var}(X_{r,j}) \le \sum_{j:f_j > 0} \mathbb{E}(X_{r,j}^2) = \sum_{j:f_j > 0} \mathbb{E}(X_{r,j}) = \frac{d}{2^r}$$

By Markov's inequality,

$$\Pr(Y_r > 0) = \Pr(Y_r \ge 1) \le \frac{\mathbb{E}(Y_r)}{1} = \frac{d}{2^r}$$

By Chebyshev's inequality,

$$\Pr(Y_r=0) \leq \Pr(|Y_r - \mathbb{E}(Y_r)| \geq d/2^r) \leq \frac{\operatorname{var}(Y_r)}{(d/2^r)^2} \leq \frac{2^r}{d}$$

For simplicity, assume n is a power of 2.

$$\mathbb{E}(Y_r) = \sum_{j:f_j > 0} \mathbb{E}(X_{r,j}) = \frac{d}{2^r}$$
$$\operatorname{var}(Y_r) = \sum_{j:f_j > 0} \operatorname{var}(X_{r,j}) \le \sum_{j:f_j > 0} \mathbb{E}(X_{r,j}^2) = \sum_{j:f_j > 0} \mathbb{E}(X_{r,j}) = \frac{d}{2^r}$$

By Markov's inequality,

$$\Pr(Y_r > 0) = \Pr(Y_r \ge 1) \le \frac{\mathbb{E}(Y_r)}{1} = \frac{d}{2^r}$$

By Chebyshev's inequality,

$$\Pr(Y_r=0) \leq \Pr(|Y_r - \mathbb{E}(Y_r)| \geq d/2^r) \leq \frac{\operatorname{var}(Y_r)}{(d/2^r)^2} \leq \frac{2^r}{d}$$

The tidemark algorithm - analysis III Recall: $\hat{d} = 2^{T+1/2}$, $\Pr(Y_r > 0) \le \frac{d}{2^r}$ and $\Pr(Y_r = 0) \le \frac{2^r}{d}$.

Recall: $\hat{d} = 2^{T+1/2}$, $\Pr(Y_r > 0) \le \frac{d}{2^r}$ and $\Pr(Y_r = 0) \le \frac{2^r}{d}$.

• Let *a* be the smallest integer such that $2^{a+1/2} \ge 3d$. We have,

$$\Pr(\hat{d} \ge 3d) = \Pr(T \ge a)$$

Recall: $\hat{d} = 2^{T+1/2}$, $\Pr(Y_r > 0) \le \frac{d}{2^r}$ and $\Pr(Y_r = 0) \le \frac{2^r}{d}$.

• Let *a* be the smallest integer such that $2^{a+1/2} \ge 3d$. We have,

$$\Pr(\hat{d} \ge 3d) = \Pr(T \ge a) = \Pr(Y_a > 0) \le \frac{d}{2^a} \le \frac{\sqrt{2}}{3}$$

Recall: $\hat{d} = 2^{T+1/2}$, $\Pr(Y_r > 0) \le \frac{d}{2^r}$ and $\Pr(Y_r = 0) \le \frac{2^r}{d}$.

• Let *a* be the smallest integer such that $2^{a+1/2} \ge 3d$. We have,



Recall: $\hat{d} = 2^{T+1/2}$, $\Pr(Y_r > 0) \le \frac{d}{2^r}$ and $\Pr(Y_r = 0) \le \frac{2^r}{d}$.

• Let *a* be the smallest integer such that $2^{a+1/2} \ge 3d$. We have,

$$\Pr(\hat{d} \ge 3d) = \Pr(T \ge a) = \Pr(Y_a > 0) \le \frac{d}{2^a} \le \frac{\sqrt{2}}{3}$$

This gives us the probability that our estimate is too large. We now bound the probability it is too small.

Recall: $\hat{d} = 2^{T+1/2}$, $\Pr(Y_r > 0) \le \frac{d}{2^r}$ and $\Pr(Y_r = 0) \le \frac{2^r}{d}$.

• Let *a* be the smallest integer such that $2^{a+1/2} \ge 3d$. We have,

$$\Pr(\hat{d} \ge 3d) = \Pr(T \ge a) = \Pr(Y_a > 0) \le \frac{d}{2^a} \le \frac{\sqrt{2}}{3}$$

This gives us the probability that our estimate is too large. We now bound the probability it is too small.

For the probability that our estimate is too small let b be the largest integer such that 2^{b+1/2} ≤ d/3.

$$\Pr(\hat{d} \le d/3) = \Pr(T \le b) = \Pr(Y_{b+1} = 0) \le \frac{2^{b+1}}{d} \le \frac{\sqrt{2}}{3}$$

Recall: $\hat{d} = 2^{T+1/2}$, $\Pr(Y_r > 0) \le \frac{d}{2^r}$ and $\Pr(Y_r = 0) \le \frac{2^r}{d}$.

• Let *a* be the smallest integer such that $2^{a+1/2} \ge 3d$. We have,

$$\Pr(\hat{d} \ge 3d) = \Pr(T \ge a) = \Pr(Y_a > 0) \le \frac{d}{2^a} \le \frac{\sqrt{2}}{3}$$

This gives us the probability that our estimate is too large. We now bound the probability it is too small.

For the probability that our estimate is too small let b be the largest integer such that 2^{b+1/2} ≤ d/3.

$$\Pr(\hat{d} \le d/3) = \Pr(T \le b) = \Pr(Y_{b+1} = 0) \le \frac{2^{b+1}}{d} \le \frac{\sqrt{2}}{3}$$

The bounds are not ideal in two ways.

- 1. We can't get an arbitrarily close approximation (yet).
- 2. The failure probably is high. $\sqrt{2}/3 \approx 0.47!$

 $\rm TIDEMARK\ space/time$

```
initialise

Choose random h: [n] \rightarrow [n]

Set z = 0

TIDEMARK(a_i)

if ZEROS(h(a_i)) > z

set z = ZEROS(h(a_i))

OUTPUT 2^{z+\frac{1}{2}}
```

```
initialise

Choose random h: [n] \rightarrow [n]

Set z = 0

TIDEMARK(a_i)

if ZEROS(h(a_i)) > z

set z = ZEROS(h(a_i))

OUTPUT 2^{z+\frac{1}{2}}
```

• One-pass and O(m) time. \checkmark

$T {\rm IDEMARK} \ \text{space} / \text{time}$

```
initialise
Choose random h: [n] \rightarrow [n]
Set z = 0
TIDEMARK(a_i)
if ZEROS(h(a_i)) > z
set z = ZEROS(h(a_i))
OUTPUT 2^{z+\frac{1}{2}}
```

- One-pass and O(m) time. \checkmark
- $O(\log \log n)$ bits of space. \checkmark

```
initialise

Choose random h: [n] \rightarrow [n]

Set z = 0

TIDEMARK(a_i)

if ZEROS(h(a_i)) > z

set z = ZEROS(h(a_i))

OUTPUT 2^{z+\frac{1}{2}}
```

- One-pass and O(m) time. \checkmark
- $O(\log \log n)$ bits of space. \checkmark

Why is it $O(\log \log n)$ bits of space?



- One-pass and O(m) time. \checkmark
- $O(\log \log n)$ bits of space. \checkmark

Why is it $O(\log \log n)$ bits of space?

• We store one value of $ZEROS(h(a_i))$ where $h(a_i) \in [n]$.



- One-pass and O(m) time. \checkmark
- $O(\log \log n)$ bits of space. \checkmark

Why is it $O(\log \log n)$ bits of space?

- ▶ We store one value of $ZEROS(h(a_i))$ where $h(a_i) \in [n]$.
- ▶ The maximum value of $\operatorname{ZEROS}(h(a_i))$ is $\lfloor \log_2 n \rfloor \in O(\log n)$.



- One-pass and O(m) time. \checkmark
- $O(\log \log n)$ bits of space. \checkmark

Why is it $O(\log \log n)$ bits of space?

- ▶ We store one value of $ZEROS(h(a_i))$ where $h(a_i) \in [n]$.
- ▶ The maximum value of $\operatorname{ZEROS}(h(a_i))$ is $\lfloor \log_2 n \rfloor \in O(\log n)$.
- Therefore only O(log log n) bits are needed to represent the biggest possible value.

Method: For each value in the input, run k independent copies of TIDEMARK in parallel and output the median.

Method: For each value in the input, run k independent copies of TIDEMARK in parallel and output the median.

Recall: $\Pr(\hat{d} \ge 3d) \le \frac{\sqrt{2}}{3}$

Method: For each value in the input, run k independent copies of TIDEMARK in parallel and output the median.

Recall:
$$\Pr(\hat{d} \ge 3d) \le \frac{\sqrt{2}}{3}$$

If the median is greater than 3d, then at least k/2 values are at least 3d.

Method: For each value in the input, run k independent copies of TIDEMARK in parallel and output the median.

Recall:
$$\Pr(\hat{d} \ge 3d) \le \frac{\sqrt{2}}{3}$$

If the median is greater than 3d, then at least k/2 values are at least 3d.

Define $X_i = 1$ if $\hat{d} \ge 3d$, 0 otherwise. $X = \sum_{i=1}^k X_i$, $\mu \le \sqrt{2}k/3$.

Let
$$\delta = 3/(2\sqrt{2}) - 1 \approx 0.06$$
 so $(1 + \delta)\mu = k/2$.

$$\Pr[X \ge (1 + \delta)\mu] \le \exp(-\delta^2 \mu/3)$$
(Chernoff bound)

Method: For each value in the input, run k independent copies of TIDEMARK in parallel and output the median.

Recall:
$$\Pr(\hat{d} \ge 3d) \le \frac{\sqrt{2}}{3}$$

If the median is greater than 3d, then at least k/2 values are at least 3d.

Define $X_i = 1$ if $\hat{d} \ge 3d$, 0 otherwise. $X = \sum_{i=1}^k X_i$, $\mu \le \sqrt{2}k/3$.

Let
$$\delta = 3/(2\sqrt{2}) - 1 \approx 0.06$$
 so $(1 + \delta)\mu = k/2$.

$$\Pr[X \ge (1 + \delta)\mu] \le \exp(-\delta^2 \mu/3)$$
(Chernoff bound)

In other words, as k grows the probability that the median is at least 3d decreases exponentially.

The median trick - lower bound

Method: For each value in the input, run k independent copies of TIDEMARK in parallel and output the median.

Recall:
$$\Pr(\hat{d} \le d/3) \le \frac{\sqrt{2}}{3}$$

▶ If the median is less than $\frac{d}{3}$, then at least k/2 values are at most $\frac{d}{3}$.

Define $X_i = 1$ if $\hat{d} \leq d/3$, 0 otherwise, $X = \sum_{i=1}^k X_i$, $\mu \leq \sqrt{2}k/3$.

Let
$$\delta = 3/(2\sqrt{2}) - 1 \approx 0.06$$
 so $(1 + \delta)\mu = k/2$.

$$\Pr[X \ge (1 + \delta)\mu] \le \exp(-\delta^2 \mu/3)$$
(Chernoff bound)

In other words, as k grows the probability that the median is at most d/3 decreases exponentially.

We saw two streaming algorithms for estimating the number of distinct items in a stream.

They are both one-pass algorithms and run in O(m) time.

We saw two streaming algorithms for estimating the number of distinct items in a stream.

They are both one-pass algorithms and run in O(m) time.

They both hash the incoming elements to [n].

▶ SIMPLE-COUNT stores the minimum of the hashed values.

We saw two streaming algorithms for estimating the number of distinct items in a stream.

They are both one-pass algorithms and run in O(m) time.

- ▶ SIMPLE-COUNT stores the minimum of the hashed values.
- ▶ SIMPLE-COUNT uses *O*(log *n*) bits of space.

We saw two streaming algorithms for estimating the number of distinct items in a stream.

They are both one-pass algorithms and run in O(m) time.

- ▶ SIMPLE-COUNT stores the minimum of the hashed values.
- ▶ SIMPLE-COUNT uses *O*(log *n*) bits of space.
- It gives us $\Pr(d \le \hat{d}/3) \le 1/3$ and $\Pr(d > 3\hat{d}) \le 1/3$.

We saw two streaming algorithms for estimating the number of distinct items in a stream.

They are both one-pass algorithms and run in O(m) time.

- ▶ SIMPLE-COUNT stores the minimum of the hashed values.
- ▶ SIMPLE-COUNT uses *O*(log *n*) bits of space.
- It gives us $\Pr(d \le \hat{d}/3) \le 1/3$ and $\Pr(d > 3\hat{d}) \le 1/3$.
- TIDEMARK stores the maximum number of trailing zeros in the hashed values.

We saw two streaming algorithms for estimating the number of distinct items in a stream.

They are both one-pass algorithms and run in O(m) time.

- ▶ SIMPLE-COUNT stores the minimum of the hashed values.
- ▶ SIMPLE-COUNT uses *O*(log *n*) bits of space.
- It gives us $\Pr(d \le \hat{d}/3) \le 1/3$ and $\Pr(d > 3\hat{d}) \le 1/3$.
- TIDEMARK stores the maximum number of trailing zeros in the hashed values.
- ▶ TIDEMARK uses O(log log n) bits of space.

We saw two streaming algorithms for estimating the number of distinct items in a stream.

They are both one-pass algorithms and run in O(m) time.

- ▶ SIMPLE-COUNT stores the minimum of the hashed values.
- ▶ SIMPLE-COUNT uses *O*(log *n*) bits of space.
- It gives us $\Pr(d \le \hat{d}/3) \le 1/3$ and $\Pr(d > 3\hat{d}) \le 1/3$.
- TIDEMARK stores the maximum number of trailing zeros in the hashed values.
- ▶ TIDEMARK uses O(log log n) bits of space.
- It gives us $\Pr(\hat{d} \le d/3) \le \frac{\sqrt{2}}{3}$ and $\Pr(\hat{d} \ge 3d) \le \frac{\sqrt{2}}{3}$

We saw two streaming algorithms for estimating the number of distinct items in a stream.

They are both one-pass algorithms and run in O(m) time.

- ▶ SIMPLE-COUNT stores the minimum of the hashed values.
- ▶ SIMPLE-COUNT uses *O*(log *n*) bits of space.
- It gives us $\Pr(d \le \hat{d}/3) \le 1/3$ and $\Pr(d > 3\hat{d}) \le 1/3$.
- TIDEMARK stores the maximum number of trailing zeros in the hashed values.
- ▶ TIDEMARK uses O(log log n) bits of space.
- It gives us $\Pr(\hat{d} \le d/3) \le \frac{\sqrt{2}}{3}$ and $\Pr(\hat{d} \ge 3d) \le \frac{\sqrt{2}}{3}$
- By performing k parallel iterations the probability of being outside the range [d/3, 3d] can be made exponentially small.