# Approximation Techniques

Supplemental Reading
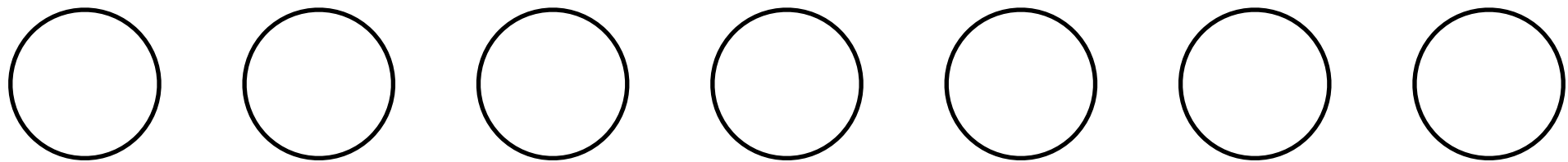(papers posted on Piazza)

1

# Announcements

- First P3 Time-Trial-Submissions Tonight.

- P2 Grades Out By Friday.

- Homework 7 due Monday.

2

# Review:WINDOW

```
SELECT L.state, T.month,
       AVG(S.sales) OVER W as movavg
FROM   Sales S, Times T, Locations L
WHERE  S.timeid = T.timeid
  AND  S.locid = L.locid
WINDOW W AS (
    PARTITION BY L.state
    ORDER BY T.month
    RANGE BETWEEN INTERVAL '1' MONTH PRECEDING
          AND INTERVAL '1' MONTH FOLLOWING
)
```
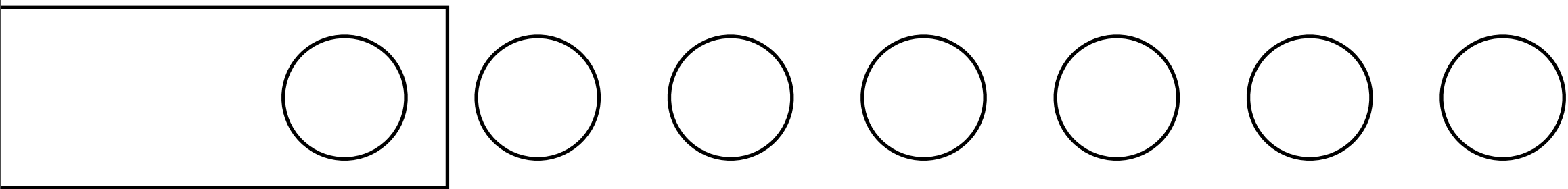
# Review: WINDOW

## Windowed SUM, Window Size 3
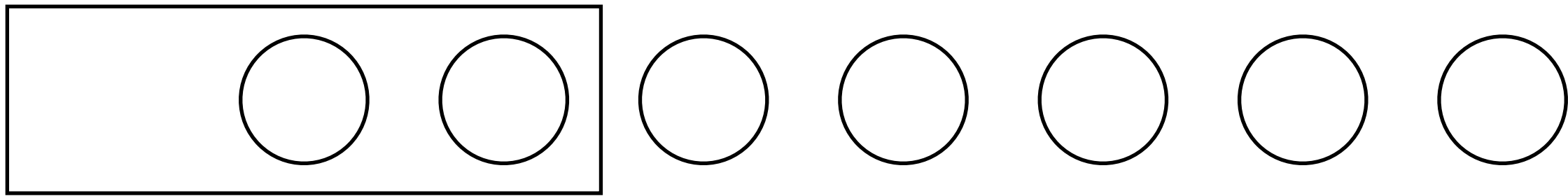
○ ○ ○ ○ ○ ○ ○

4

# Review: WINDOW

## Windowed SUM, Window Size 3

SUM₁

# Review: WINDOW

## Windowed SUM, Window Size 3

SUM$_1$ SUM$_2$

4

# Review: WINDOW

## Windowed SUM, Window Size 3

○ ○ ○ ○ ○ ○ ○

$SUM_1$ $SUM_2$ $SUM_3$ $SUM_4$ $SUM_5$ $SUM_6$ $SUM_7$ …

4

# Review: WINDOW

```
SELECT L.state, T.month,
       AVG(S.sales) OVER W as movavg
FROM   Sales S, Times T, Locations L
WHERE  S.timeid = T.timeid
  AND  S.locid = L.locid
WINDOW W AS (
   PARTITION BY L.state
   ORDER BY T.month
   RANGE BETWEEN INTERVAL '1' MONTH PRECEDING
         AND INTERVAL '1' MONTH FOLLOWING
)
```

Partition By is like Group By

Required: Define a sort order

Required: Define the size of window
(need not be a fixed # of tuples)

5

# Review: Data Warehousing

- Data warehouses store massive datasets

- Workload involves…

  - Frequent, low-latency reads.

  - Lots of aggregation/summarization.

6

# Review: Data Warehousing

- Data warehouses store massive datasets

- Workload involves…

  - Frequent, low-latency reads.

  - Lots of aggregation/summarization.

  - … and may not require precise answers.
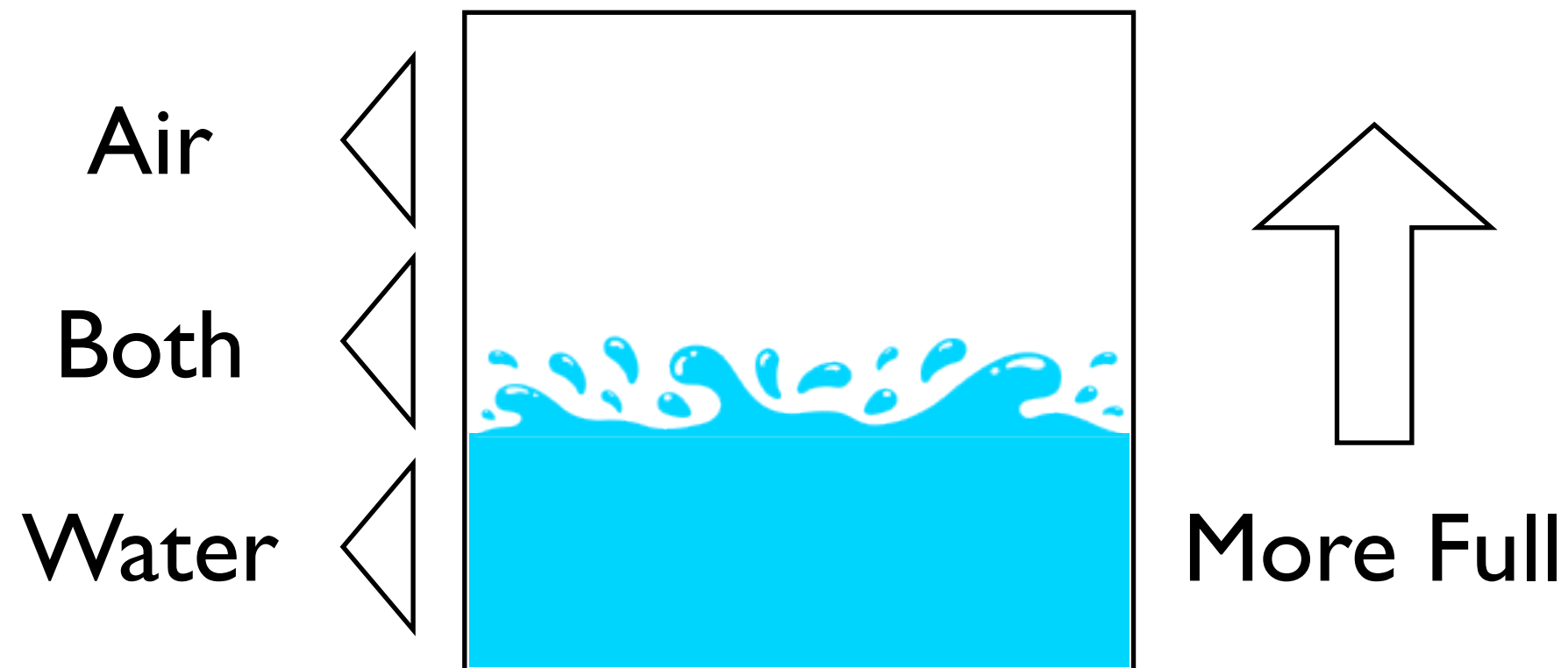
6

# Query Approximation

- Summarize the data with a **Sketching** Algorithm.

    - Bloom Filters (Set Containment)

    - Flajolet & Martin Count-Distinct Sketch

    - Count Sketch (Frequent Items/Top-K)

- Sample the data and estimate the error

    - … or better yet, keep generating samples!

        - Online Aggregation & Ripple Joins

7

# Count-Distinct Sketches

- Count-Distinct is Holistic (all tuples needed)

  - Naive External Algorithm: Sort/Eliminate

- Flajolet & Martin Sketch:

  - Summarize dataset in a bit vector.

  - Bit vector 'fills up' as more values added.

  - Hashes eliminate duplicate contributions.

8

# Count-Distinct Sketches



Air

Both

Water

More Full

image credit: openclipart.org

Wednesday, April 10, 13

# Count-Distinct Sketches

0s      Both           1s

0001  0100  0111  1111

'Fill' Boundaries Approximate # of Distinct Items in Set

# The ρ function

hash($\Diamond$) = `00010010100`

P[ ρ($\Diamond$) = 0 ] = ?

P[ ρ($\Diamond$) = 1 ] = ?

P[ ρ($\Diamond$) = k ] = ?

11

# The ρ function

hash($\diamondsuit$) = 00010010100

$\uparrow$

Smallest Position with a Non-Zero Bit
(or |bit vector|)

ρ($\diamondsuit$) = 2

P[ ρ($\diamondsuit$) = 0 ] = ?

P[ ρ($\diamondsuit$) = 1 ] = ?

P[ ρ($\diamondsuit$) = k ] = ?

11

# The ρ function

$$\text{hash}(\diamondsuit) = \texttt{00010010100}$$

$$\text{sketch}(\diamondsuit) = 2^{\rho(\diamondsuit)}$$

$$= \texttt{00000000100}$$

$$\text{sketch}(\diamondsuit_1, \diamondsuit_2, \dots) = \text{sketch}(\diamondsuit_1) \vee \text{sketch}(\diamondsuit_2) \vee \dots$$

12

# Count-Distinct Sketches

Each item $\diamondsuit$ has a 1/2 chance of $\rho(\diamondsuit) = 0$

What is the probability that $\rho(\diamondsuit) \neq 0$ for all N items?

# Count-Distinct Sketches

Each item $\diamondsuit$ has a 1/2 chance of $\rho(\diamondsuit) = 0$

What is the probability that $\rho(\diamondsuit) \neq 0$ for all N items?

Each item $\diamondsuit$ has a $1/2^k$ chance of $\rho(\diamondsuit) = k$

What is the probability that $\rho(\diamondsuit) \neq k$ for all N items?

13

# Count-Distinct Sketches

Given a sketch bit vector,
let <u>R</u> be the position of the lowest zero value.

```
00010111111
```

$$E[\ R\ ] = \log_2(\ \varphi\ *\ |Set|\ )$$
$$\varphi = 0.77351$$

# Count-Distinct Sketches

Given a sketch bit vector,
let <u>R</u> be the position of the lowest zero value.

```
00010111111
```
↑

R=6

$$E[ R ] = \log_2( \varphi * |Set| )$$

$$\varphi = 0.77351$$

14

# Count-Distinct Sketches

$$\varphi = 0.77351$$

$$E[\ R\ ] = \log_2(\ \varphi\ *\ |Set|\ )$$

$$2^{R/\varphi} = E[\ |Set|\ ]$$

$$2^{6/\varphi} = 60$$

0001011111 Summarizes a set with
60 distinct elements

# Count-Distinct Sketches

- Problem: Estimate has a high variance.

  - Solution: Multiple sketches in parallel.

    - Use average or median of all estimates.

- Question: How does this algorithm count only unique values (count <u>distinct</u>)?
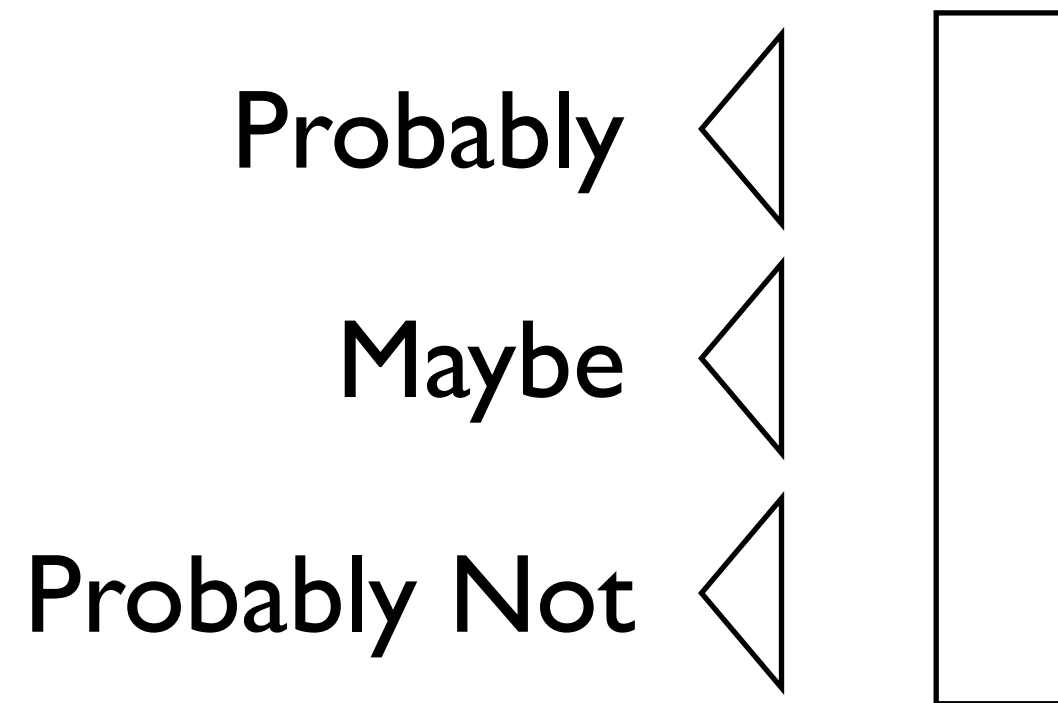
- Question: How big is the bit vector?

16

# Count Sketch

- Top-K-Count

  - Compute Group-by Count Aggregate.

  - Find the K items with the highest counts.

- Top-K-Count is Holistic

  - Sketch provides some "wiggle room"

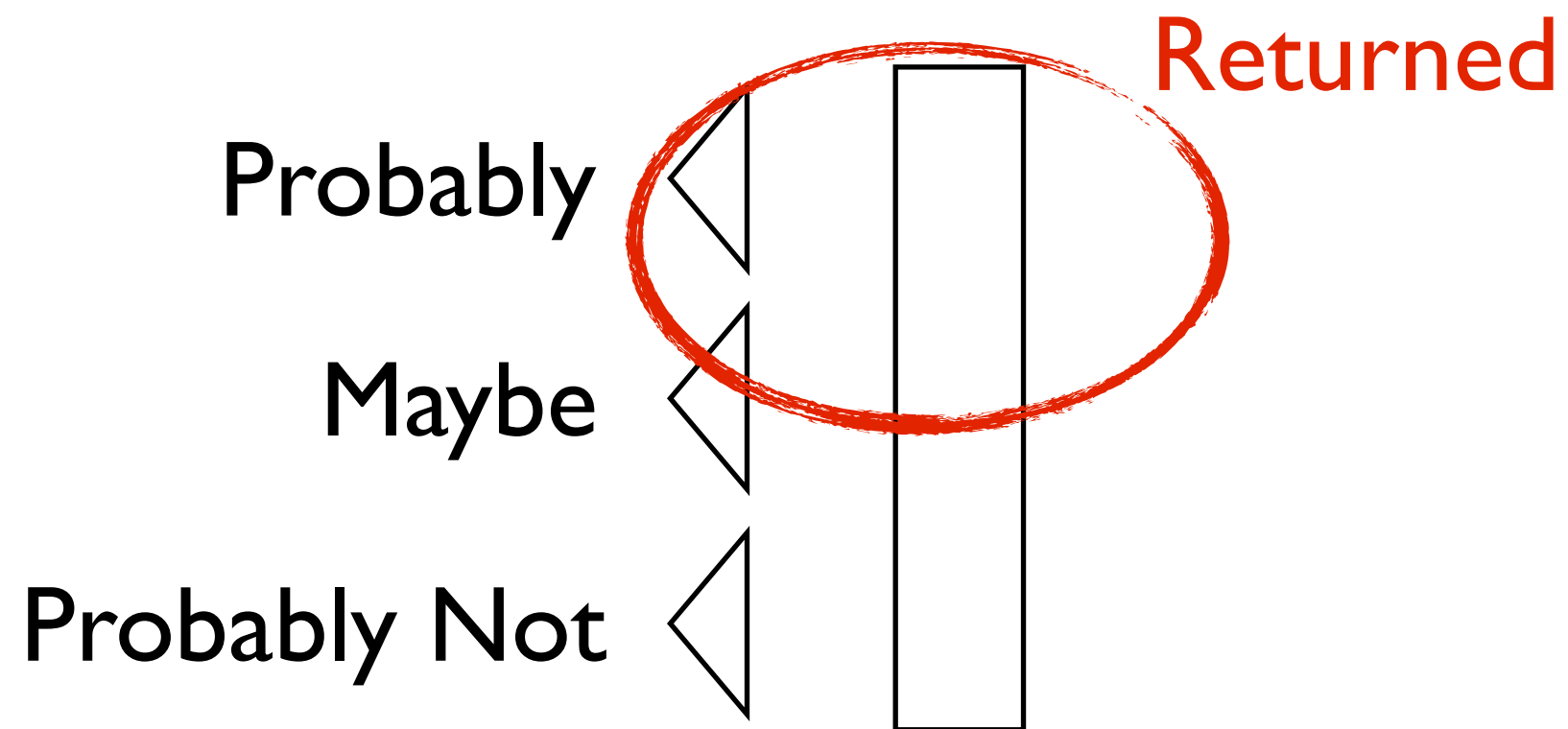    - "Borderline" entries may be excluded.

17

# Count Sketch

- Count Sketch Parameters:

  - k : The number of elements to return

  - $\epsilon$ : The 'wiggle room'

- Count Sketch Guarantees:

  - If $n_k$ is the lowest count in the top k, all objects returned have count $n_i > (1-\epsilon)\, n_k$.

  - w.h.p., all objects with $n_i > (1+\epsilon)\, n_k$ are returned.

18

# Count Sketch

Probably

Maybe

Probably Not

19

# Count Sketch



Probably

Maybe

Probably Not

Returned

# Intuition

K entries
(Tracked Explicitly)

For each tuple…

Update a count
Move the tuple up
Update the approximation

Everything
Else
(Approximated)

20

# Count Sketch

$$s(\diamond) \rightarrow \{ +1, -1 \} \ (1\text{-bit of hash}(\diamond))$$

$$\text{sketch}(\diamond_1, \diamond_2, \ldots) = \sum_i s(\diamond_i)$$

# Count Sketch

$$s(\diamondsuit) \rightarrow \{ +1, -1 \} \text{ (1-bit of hash}(\diamondsuit))$$

$$\text{sketch}(\diamondsuit_1, \diamondsuit_2, \ldots) = \sum_i s(\diamondsuit_i)$$

For a set R containing precisely N instances of $\diamondsuit$
and nothing else, what is E[ sketch(R) ]?

21

# Count Sketch

$$s(\diamondsuit) \rightarrow \{ +1, -1 \} \text{ (1-bit of hash}(\diamondsuit))$$

$$\text{sketch}(\diamondsuit_1, \diamondsuit_2, \dots) = \sum_i s(\diamondsuit_i)$$

For a set R containing precisely N instances of $\diamondsuit$
and nothing else, what is E[ sketch(R) ]?

For a set R containing an entirely <u>random</u> set of elements
what is E[ sketch(R) ]?

21

# Count Sketch

$$E[count(\lozenge)] = s(\lozenge) * sketch(R)$$

22

# Count Sketch

$$E[count(\diamond)] = s(\diamond) * sketch(R)$$

Correct (but very non-intuitive)
**Problem**: EXTREMELY high variance

22

# Count Sketch

$$E[count(\diamondsuit)] = s(\diamondsuit) * sketch(R)$$

Correct (but very non-intuitive)

**Problem**: EXTREMELY high variance

**Solution**: Use multiple sketches

22

# Count Sketch

| | | | | | |
|---|---|---|---|---|---|
| | | | | | |
| | | | | | |
| | | Keep N * M sketches | | | |
| | | | | | |
| | | | | | |

For each ◇
  For x in [0 to N)
    Update sketch (x, hash$_x$(◇) % M)

23

# Count Sketch

| | | | | | |
|---|---|---|---|---|---|
| | | | | | |
| | | | | | |
| | Keep N * M sketches | | | | |
| | | | | | |
| | | | | | |
| | | | | | |

To get approximate count($\diamondsuit$):

    For x in [0 to N)

        Approximate with sketch (x, $hash_x(\diamondsuit)$ % M)
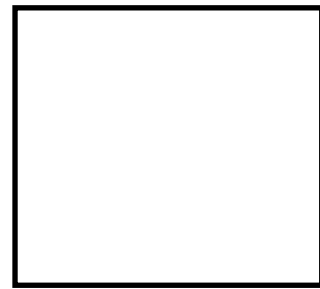
    Take <u>median</u> approximation over all x

24

# Count Sketch

Variance is still high.

But good enough for Top-K

25

# Count Sketch

K entries
(Tracked Explicitly)

Everything
Else
(Count Sketch)

For each tuple T:

Tuple already in Top K?

Update count of T

Otherwise

E[count T] > min count of all tuples in Top K?

T replaces lowest

Update Count Sketch

26

# Sketching Algorithms

- Summarize data in a fixed amount of space.

  - Specialized for a specific aggregate.

- Provide probabilistic guarantees.

  - Use properties of how sketch is updated.

  - Use hash and idempotent sketch updates to deduplicate values from the set.

27

# Bibliography

- **Probabilistic Counting Algorithms for Data Bases**

  - **Flajolet and Martin**

- **Finding Frequent Items in Data Streams**

  - **Charikar, Chen and Farach-Colton**

- Online Aggregation

  - Hellerstein, Haas, Wang

- Ripple Joins For Online Aggregation

  - Haas, Hellerstein