

Summarizing and mining inverse distributions on data streams via dynamic inverse sampling

Presented by

Graham Cormode

`cormode@bell-labs.com`

S. Muthukrishnan

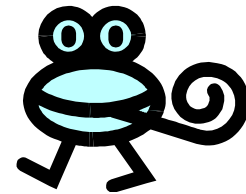
`muthu@cs.rutgers.edu`

Irina Rozenbaum

`rozenbau@paul.rutgers.edu`

Outline

- Defining and motivating the Inverse Distribution
- Queries and challenges on the Inverse Distribution
- Dynamic Inverse Sampling to draw sample from Inverse Distribution
- Experimental Study



Data Streams & DSMSs

- Numerous real world applications generate data streams:
 - IP network monitoring
 - click streams
 - Telecommunications
 - financial transactions
 - sensor networks
 - text streams at application level, etc.
- Data streams are characterized by massive data volumes of transactions and measurements at high speeds.
- Query processing is difficult on data streams:
 - We cannot store everything, and must process at line speed.
 - Exact answers to many questions are impossible without storing everything
 - We must use approximation and randomization with strong guarantees.
- Data Stream Management Systems (DSMS) summarize streams in small space (samples and sketches).

DSMS Application: IP Network Monitoring



- Needed for:
 - network traffic patterns identification
 - intrusion detection
 - reports generation, etc.



- IP traffic stream:

- Massive data volumes of transactions and measurements:
 - over 50 billion flows/day in AT&T backbone.
- Records arrive at a fast rate:
 - DDoS attacks - up to 600,000 packets/sec



- Query examples:

- heavy hitters
- change detection
- quantiles
- Histogram summaries

REF. NO.	ITEM	DATE	DESCRIPTION	AMOUNT	BALANCE
			Personal Withdrawal - 5		
484	12-MAR-2010		---Currently due Home--- 6	40.00	40.00
485	12-MAR-2010		SCM's Parking Charge Incident	386.00	386.00
486	12-MAR-2010		Balance Fwd	2.00	2.00
			Charge Materials & Supplies	2.00	2.00
			---Feeding Standard 484--- 7		
20000			Direct Staffed Loan Bal	-970.00	-970.00
20000			Direct Staffed Loan Transfer	-487.00	-487.00
			---Future due Home--- 8		
438	01-2010-2010		Short Term Loan Service Charge	6.00	6.00
477	01-2010-2010		Short Term Loan Issue	200.00	200.00

For Correspondence:	U of O Accounts Receivable
100 Williams St. Eugene, OR 97401	PO Box 3227 Eugene, OR 97431-3227
Minimum Due: 4.00	Phone: (541) 345-2177
Total Due: 2346.00	Fax: (541) 345-2177
Account Balance: 2842.00	Account Number: 123-44-5555
	Bill Date: 10-MAR-2010

Due Date:	Total Due:
01-APR-2010	\$2346.00
Account Number:	Amount Due:
123-44-5555	2346.00

U of O Accounts Receivable

100 Williams St.
Eugene, OR 97401

U of O Accounts Receivable
PO Box 3264
Portland, OR 97224-4264

Forward and Inverse Views

Consider the IP traffic on a link as packet p representing (i_p, s_p) pairs where i_p is a source IP address and s_p is a size of the packet.

Problem A.

Which IP address sent the most bytes?

That is, find i such that

$\sum_{p|i_p=i} s_p$ is maximum.

Forward distribution.

Problem B.

What is the most common volume of traffic sent by an IP address?

That is, find traffic volume W

s.t. $|\{i | W = \sum_{p|i_p=i} s_p\}|$ is maximum.

Inverse distribution.

The Inverse Distribution

If f is a discrete distribution over a large set X , then inverse distribution, $f^{-1}(i)$, gives fraction of items from X with count i .

- Inverse distribution is $f^{-1}[0\dots N]$,
 $f^{-1}(i)$ = fraction of IP addresses which sent i bytes.
 $= |\{x : f(x) = i, i^{-1} \neq 0\}| / |\{x : f(x)^{-1} \neq 0\}|$

$$F^{-1}(i) = \text{cumulative distribution of } f^{-1}$$
$$= \sum_{j > i} f^{-1}(j) \quad [\text{sum of } f^{-1}(j) \text{ above } i]$$

- Fraction of IP addresses which sent $< 1\text{KB}$ of data = $1 - F^{-1}(1024)$
- Most frequent number of bytes sent = i s.t. $f^{-1}(i)$ is greatest
- Median number of bytes sent = i s.t. $F^{-1}(i) = 0.5$

Queries on the Inverse Distribution

- Particular queries proposed in networking map onto f^{-1} ,
 - $f^{-1}(1)$ (number of flows consisting of a single packet) indicative of network abnormalities / attack [Levchenko, Paturi, Varghese 04]
 - Identify evolving attacks through shifts in Inverse Distribution [Geiger, Karamcheti, Kedem, Muthukrishnan 05]
- Better understand resource usage:
 - what is dbn. of customer traffic? How many customers < 1MB bandwidth / day? How many use 10 – 20MB per day?, etc.
→ Histograms/ quantiles on inverse distribution.
- Track most common usage patterns, for analysis / charging
 - requires heavy hitters on Inverse distribution
- Inverse distribution captures fundamental features of the distribution, has not been well-studied in data streaming.

Forward and Inverse Views on IP streams

Consider the IP traffic on a link as packet p representing (i_p, s_p) pairs where i_p is a source IP address and s_p is a size of the packet.

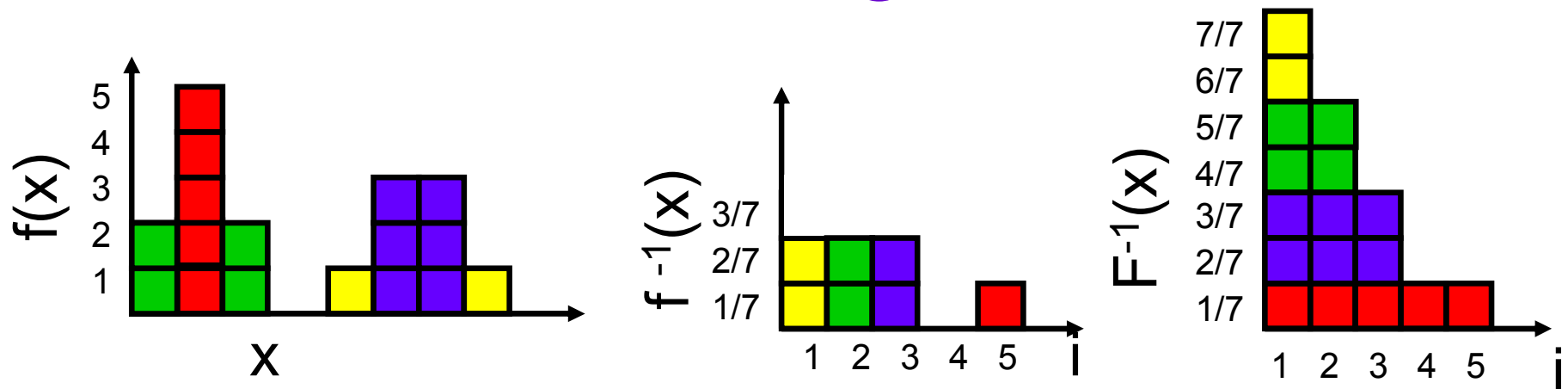
Forward distribution:

- Work on $f[0...U]$ where $f(x)$ is the number of bytes sent by IP address x .
- Each new packet (i_p, s_p) results in $f[i_p] \leftarrow f[i_p] + s_p$.
- Problems:
 - $f(i) = ?$
 - which $f(i)$ is the largest?
 - quantiles of f ?

Inverse distribution:

- Work on $f^{-1}[0...K]$
- Each new packet results in $f^{-1}[f[i_p]] \leftarrow f^{-1}[f[i_p]] - 1$ and $f^{-1}[f[i_p] + s_p] \leftarrow f^{-1}[f[i_p] + s_p] + 1$.
- Problems:
 - $f^{-1}(i) = ?$
 - which $f^{-1}(i)$ is the largest?
 - quantiles of f^{-1} ?

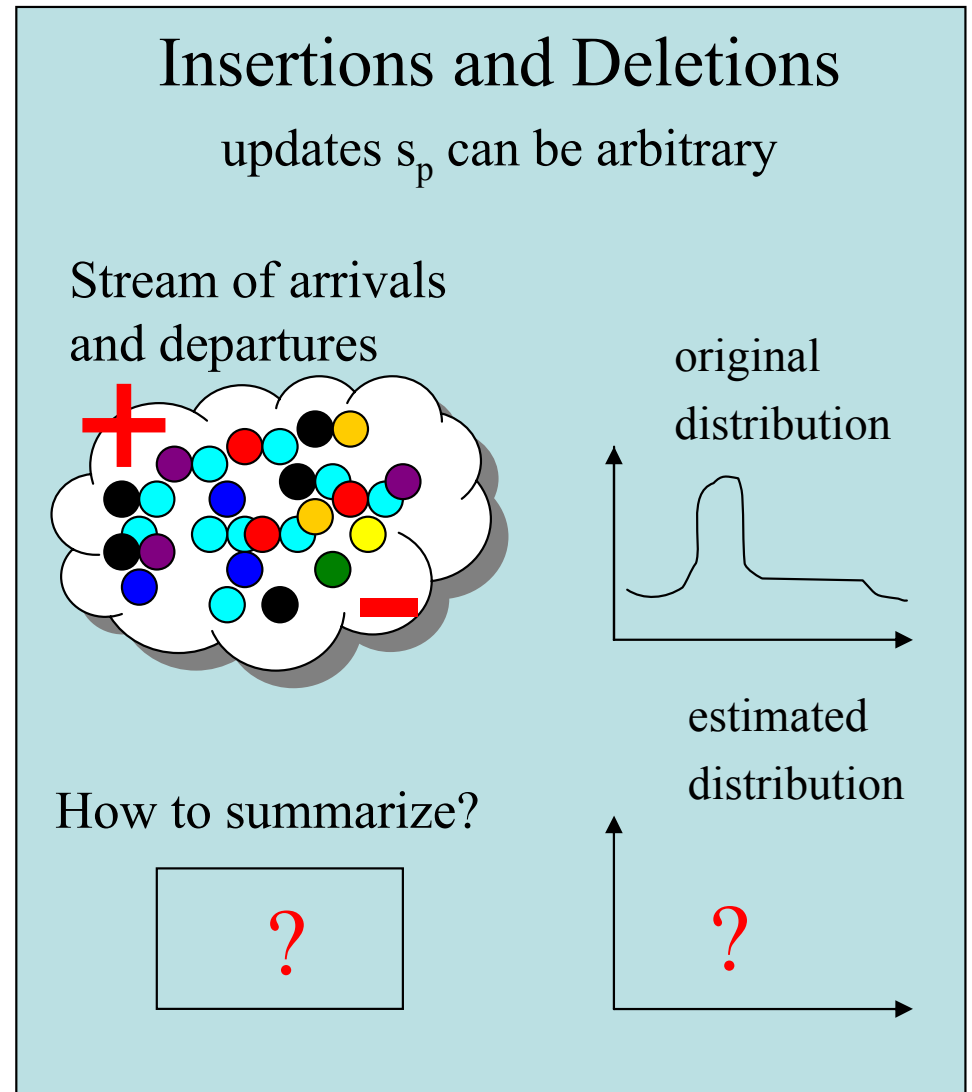
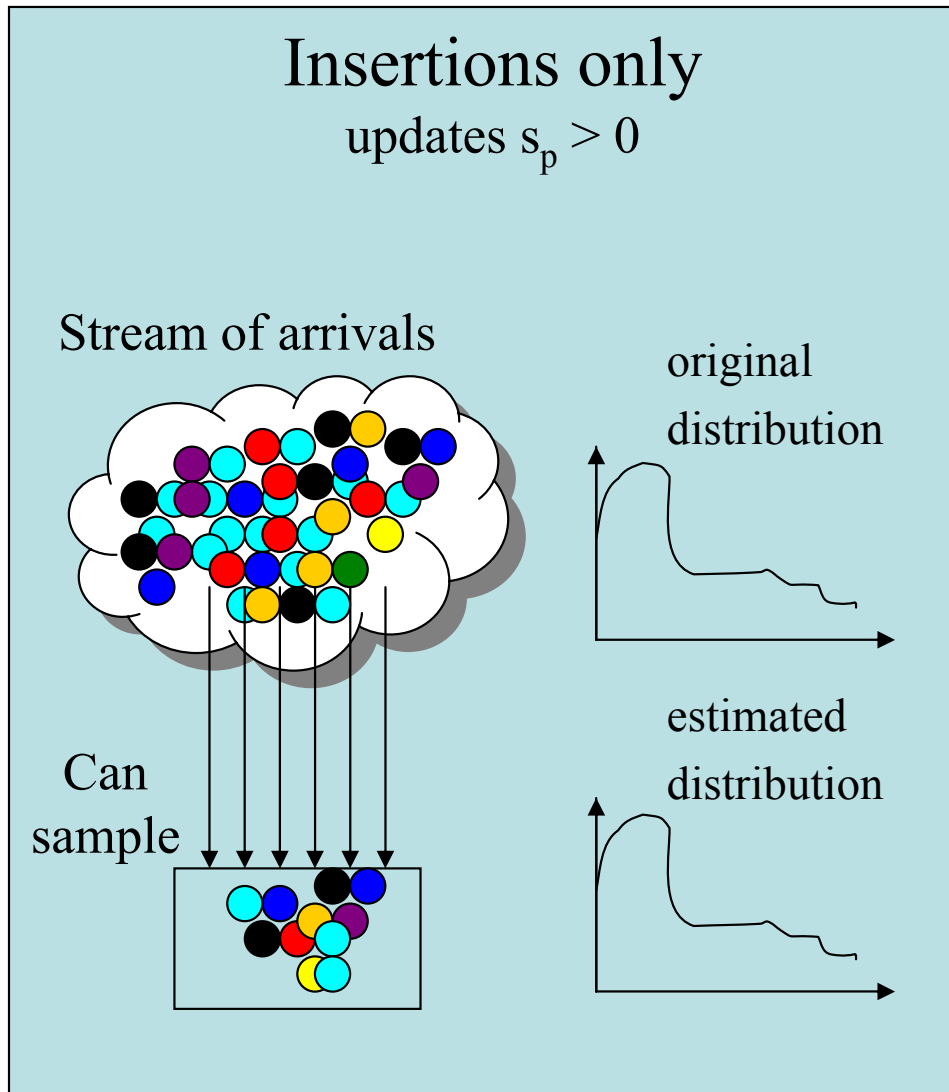
Inverse Distribution on Streams: Challenges I



- If we have full space, it is easy to go between forward and inverse distribution.
- But in small space it is much more difficult, and existing methods in small space don't apply.
- Find $f(192.168.1.1)$ in small space, with query give *a priori* – easy: just count how many times the address is seen.
- Find $f^{-1}(1024)$ – is provably hard (can't find exactly how many IP addresses sent 1KB of data without keeping full space).

Inverse Distribution on Streams: Challenges II, deletions

How to maintain summary in presence of insertions and deletions?



Our Approach: Dynamic Inverse Sampling

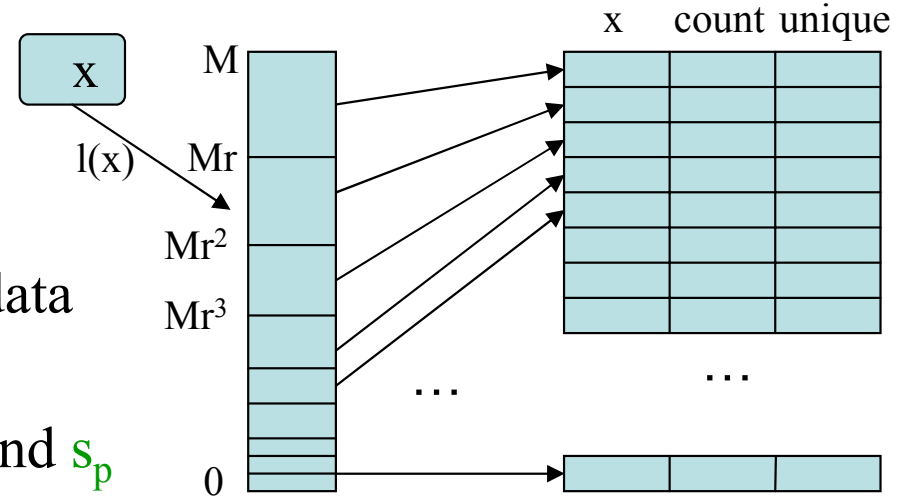
- Many queries on the forward distribution can be answered effectively by drawing a sample.
 - Draw an \mathbf{x} so probability of picking \mathbf{x} is $f(\mathbf{x}) / \sum_y f(y)$
- Similarly, we want to draw a sample from the inverse distribution in the centralized setting.
 - draw (\mathbf{i}, \mathbf{x}) s.t. $f(\mathbf{x})=\mathbf{i}$, $\mathbf{i} \neq 0$ so probability of picking \mathbf{i} is $f^{-1}(\mathbf{i}) / \sum_j f^{-1}(j)$ and probability of picking \mathbf{x} is uniform.
- Drawing from forward distribution is “easy”: just uniformly decide to sample each new item (IP address, size) seen
- Drawing from inverse distribution is more difficult, since probability of drawing $(\mathbf{i}, 1)$ should be same as $(j, 1024)$

Dynamic Inverse Sampling: Outline

- Data structure split into levels

- For each update (i_p, s_p) :

- compute hash $l(i_p)$ to a level in the data structure.
- Update counts in level $l(i_p)$ with i_p and s_p



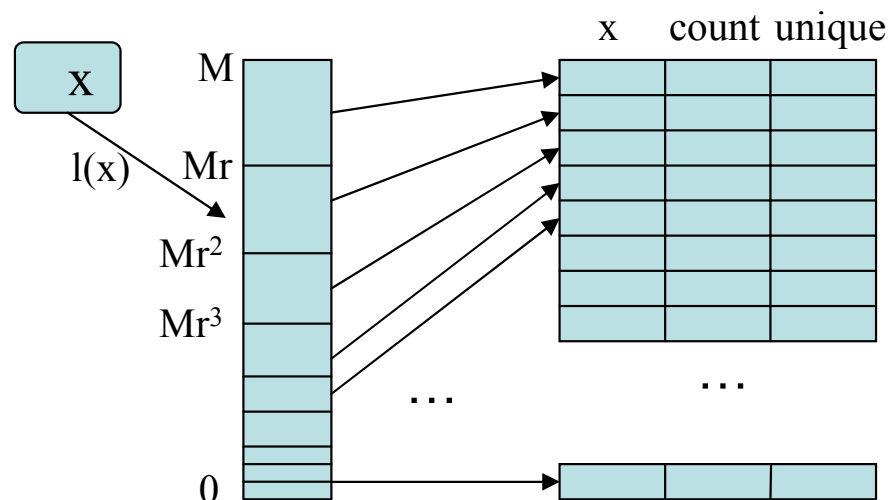
- At query time:

- probe the data structure to return $(i_p, \sum s_p)$ where i_p is sampled uniformly from all items with non-zero count
- Use the sample to answer the query on the inverse distribution.

Hashing Technique

Use hash function with exponentially decreasing distribution:
Let h be the hash function and r is an appropriate const < 1

$$\begin{aligned}\Pr[h(x) = 0] &= (1-r) \\ \Pr[h(x) = 1] &= r(1-r) \\ &\dots \\ \Pr[h(x) = l] &= r^l(1-r)\end{aligned}$$



Track the following information as updates are seen:

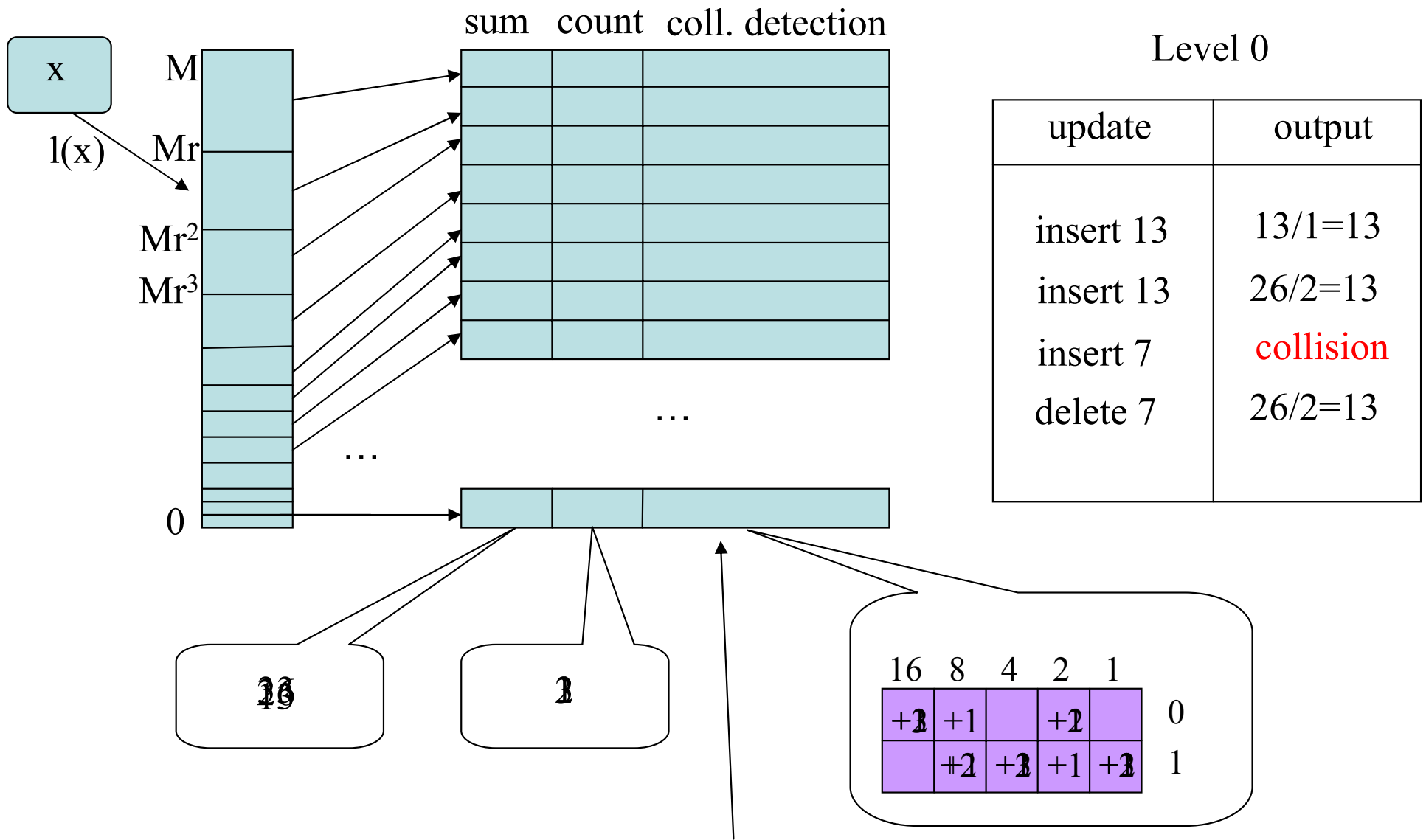
- x : Item with largest hash value seen so far
- **unique**: Is it the only distinct item seen with that hash value?
- **count**: Count of the item x

Easy to keep $(x, unique, count)$ up to date for insertions only

Challenge:

How to maintain
in presence of
deletes?

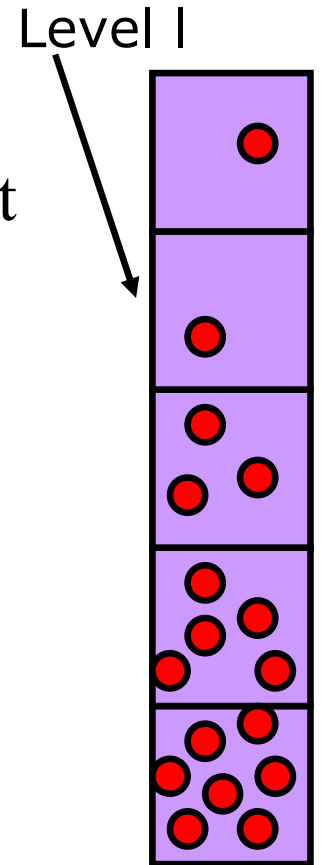
Collision Detection: inserts and deletes



Simple: Use approximate distinct element estimation routine.

Outline of Analysis

- **Analysis shows:** if there's unique item, it's chosen uniformly from set of items with non-zero count.
- Can show whatever the distribution of items, the probability of a unique item at level 1 is at least constant
- Use properties of hash function:
 - only limited, **pairwise** independence needed (easy to obtain)
- **Theorem:** With constant probability, for an arbitrary sequence of insertions and deletes, the procedure returns a uniform sample from the inverse distribution with constant probability.
- **Repeat** the process independently with different hash functions to return larger sample, with high probability.



Application to Inverse Distribution Estimates

Overall Procedure:

- Obtain the distinct sample from the inverse distribution of size s ;
- Evaluate the query on the sample and return the result.
 - Median number of bytes sent: find median from sample
 - The most common volume of traffic sent: find the most common from sample
 - What fraction of items sent i bytes: find fraction from the sample

Example:

- Median is bigger than $\frac{1}{2}$ and smaller than $\frac{1}{2}$ the values.
- Answer has some error: not $\frac{1}{2}$, but $(\frac{1}{2} \pm \epsilon)$

Theorem: If sample size $s = O(1/\epsilon^2 \log 1/\delta)$ then answer from the sample is between $(\frac{1}{2}-\epsilon)$ and $(\frac{1}{2}+\epsilon)$ with probability at least $1-\delta$.

Proof follows from application of Hoeffding's bound.

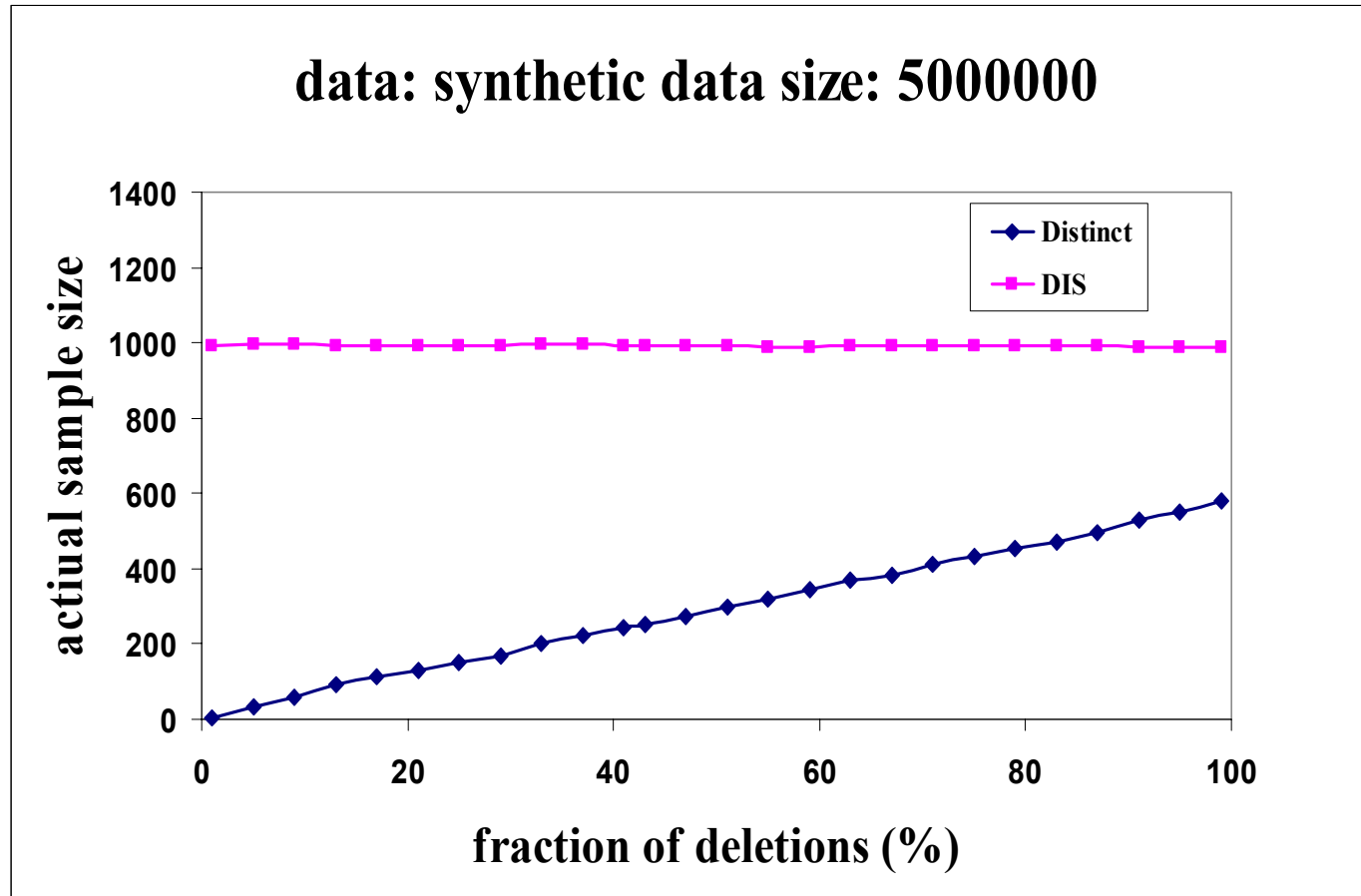
Experimental Study

Data sets:

- Large sets of network data drawn from HTTP log files from the 1998 World Cup Web Site (several million records each)
- Synthetic data set with 5 million randomly generated distinct items
- Used to build a dynamic transactions set with many insertions and deletions
- **(DIS) Dynamic Inverse Sampling algorithms** – extract at most one sample from each data structure
- **(GDIS) Greedy version of Dynamic Inverse Sampling** – greedily process every level, extract as many samples as possible from each data structure
- **(Distinct) Distinct Sampling (Gibbons VLDB 2001)** draws a sample based on a coin-tossing procedure using a pairwise-independent hash function on item values

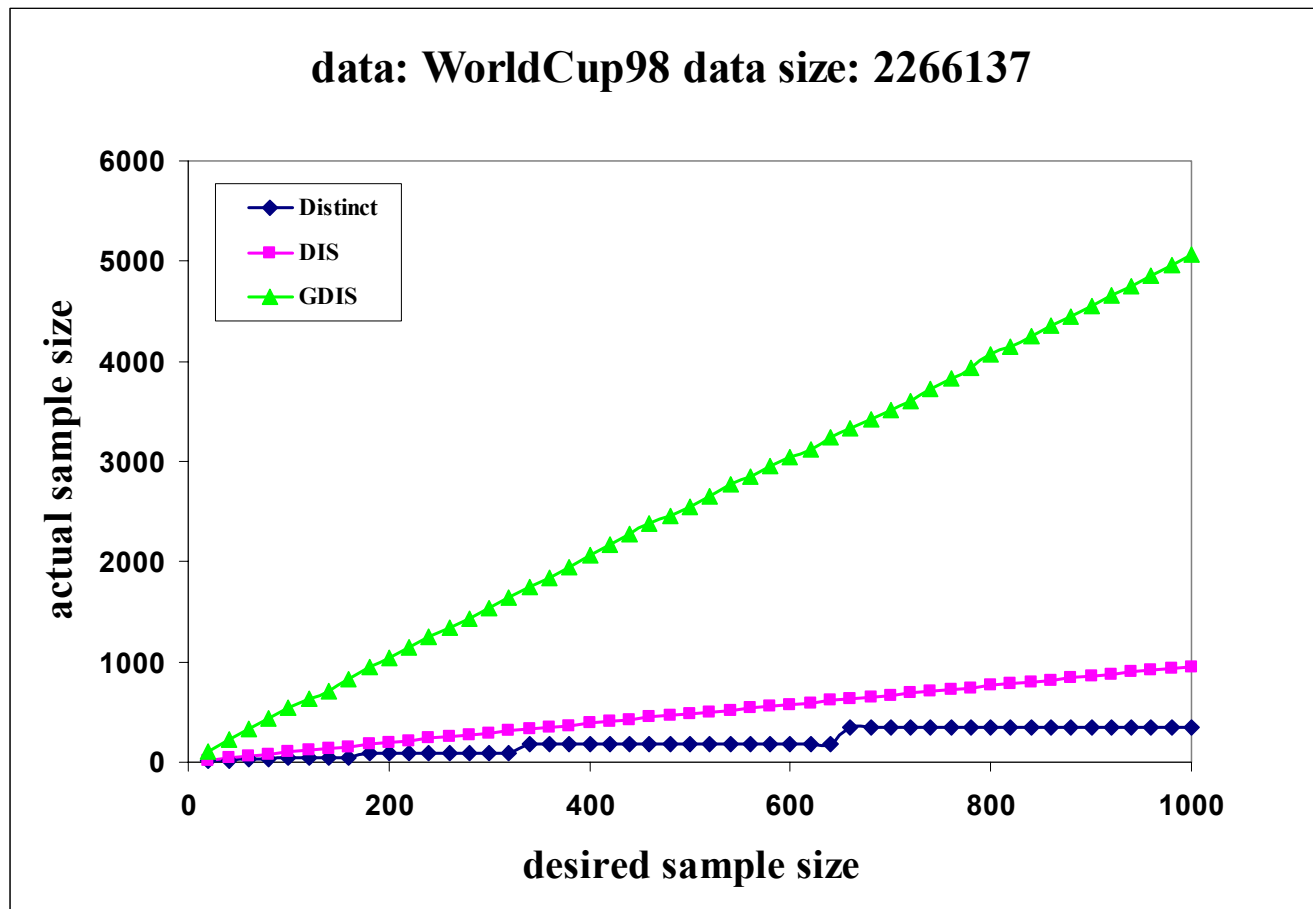
Sample Size vs. Fraction of Deletions

Desired sample size is 1000.



Returned Sample Size

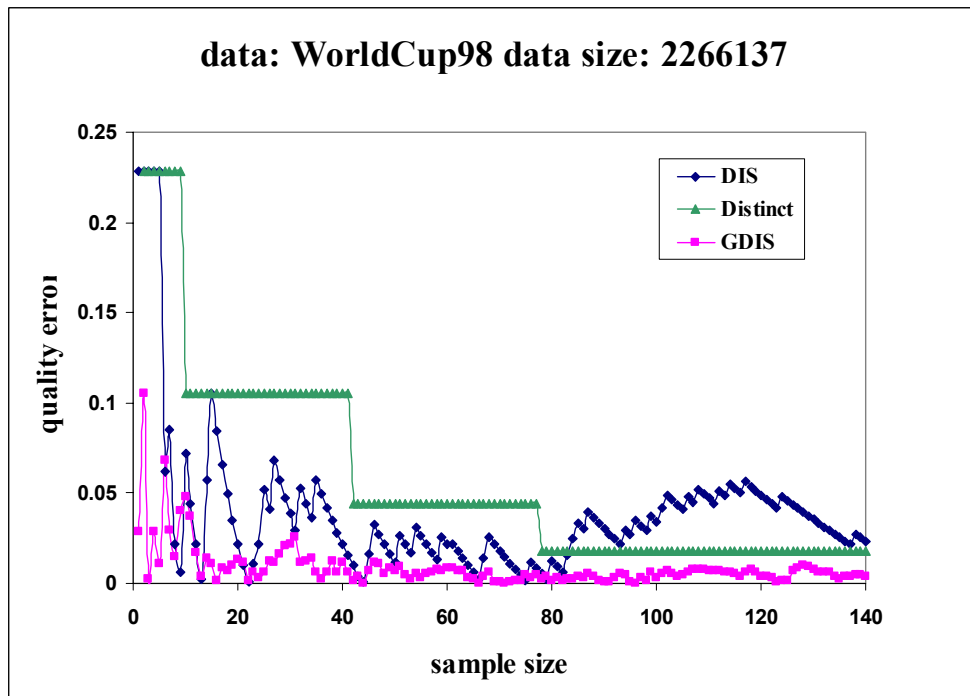
Experiments were run on the client ID attribute of the HTTP log data. 50% of the inserted records were deleted.



Sample Quality

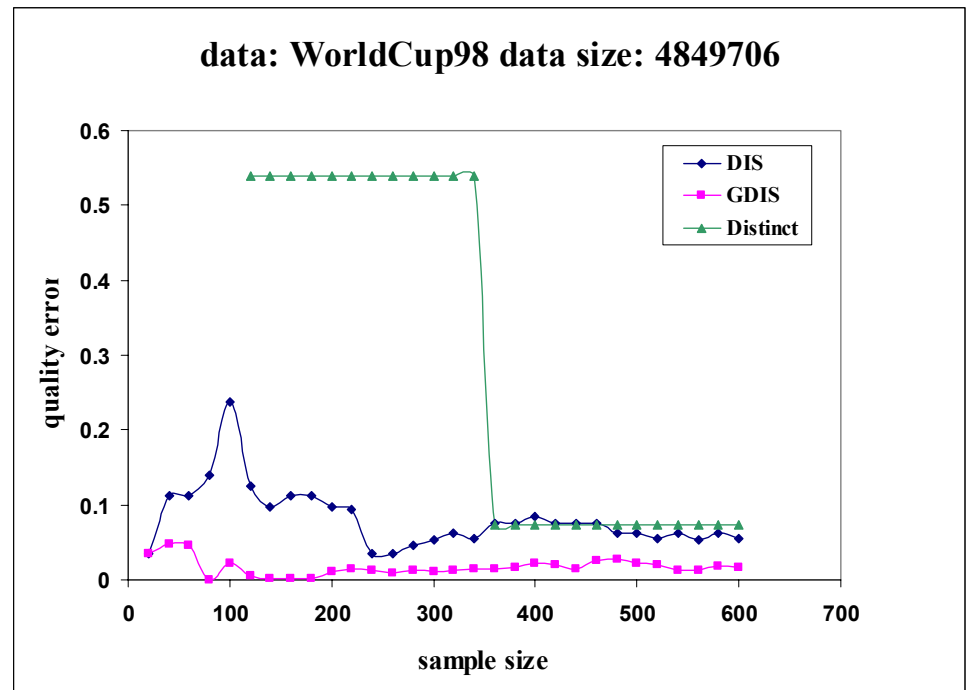
Inverse range query:

Compute the fraction of records with size greater than $i=1024$ and compare it to the exact value computed offline



Inverse quantile query:

Estimate the median of the inverse distribution using the sample and measure how far was the position of the returned item i from 0.5.



Related Work

- Distinct sampling under insert only:
 - **Gibbons**: Distinct Sampling, VLDB 2002.
 - **Datar and Muthukrishnan**: Rarity and similarity, ESA 2002.
- Distinct sampling under deletes also:
 - **Frahling, Indyk, Sohler**: Dynamic geometric streams, STOC 2005.
 - **Ganguly, Garofalakis, Rastogi**: Processing Set Expressions over Continuous Update Streams, SIGMOD 2003.
- Inverse distributions:
 - Has recently informally appeared in networking papers.

Conclusions

- We have formalized **Inverse Distributions** on data streams and introduced **Dynamic Inverse Sampling** method that draws uniform samples from the inverse distribution in presence of insertions and deletions.
- With a sample of size $O(1/\epsilon^2)$, can answer many queries on the inverse distribution (including point and range queries, heavy hitters, quantiles) up to additive approximation of ϵ .
- Experimental study shows that proposed methods can work at high rates and answer queries with high accuracy
- **Future work:**
 - Incorporate in data stream systems
 - Can we also sample from forward dbn under inserts and deletes?

