

CS60021: Scalable Data Mining

Stream Mining

Sourangshu Bhattacharya

Filtering Data Streams

Filtering Data Streams

- Each element of data stream is a tuple
- Given a list of keys S
- **Determine which tuples of stream are in S**
- **Obvious solution: Hash table**
 - But suppose we **do not have enough memory** to store all of S in a hash table
 - E.g., we might be processing millions of filters on the same stream

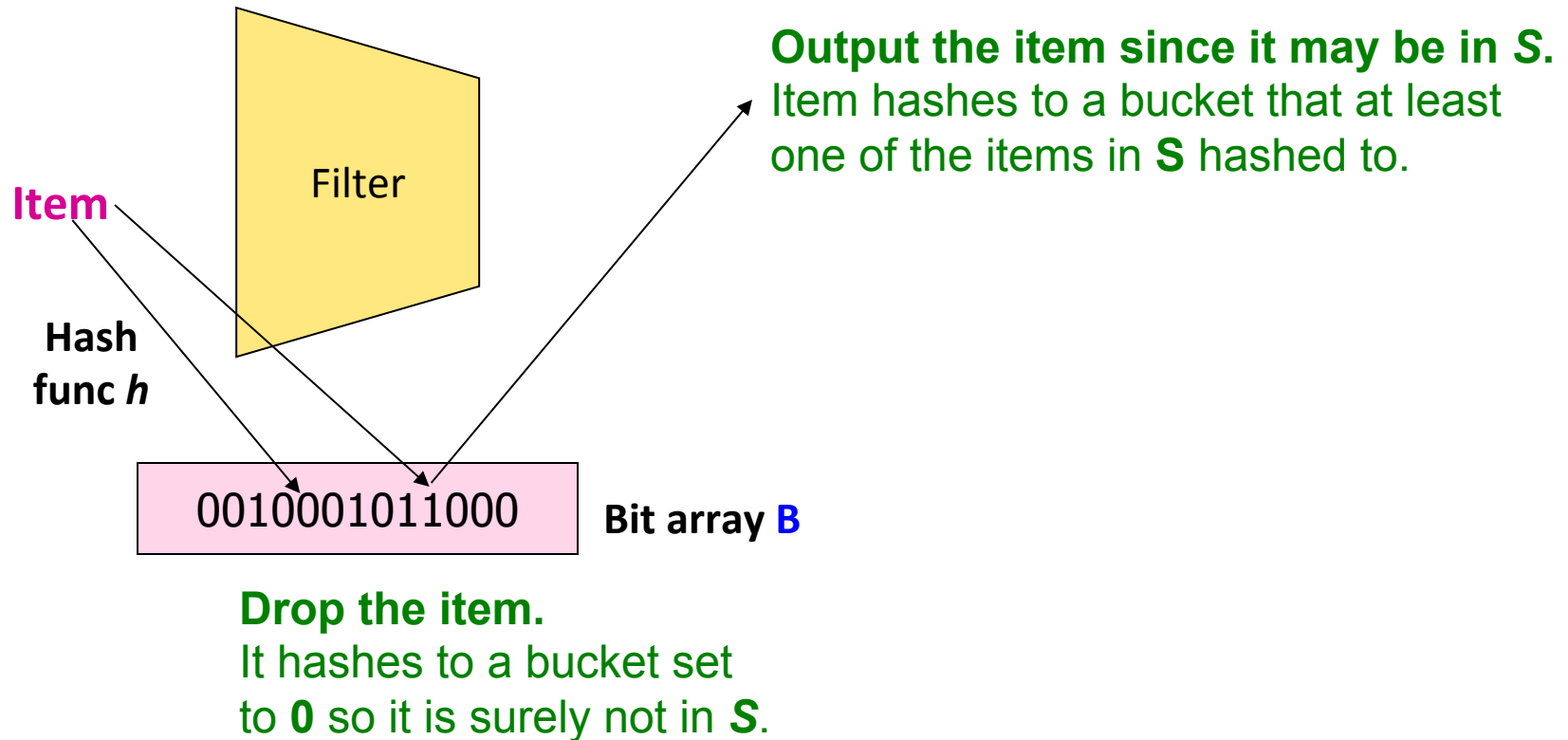
Applications

- **Example: Email spam filtering**
 - We know 1 billion “good” email addresses
 - If an email comes from one of these, it is **NOT** spam
- **Publish-subscribe systems**
 - You are collecting lots of messages (news articles)
 - People express interest in certain sets of keywords
 - Determine whether each message matches user’s interest

First Cut Solution (1)

- Given a set of keys S that we want to filter
- Create a bit array B of n bits, initially all 0s
- Choose a hash function h with range $[0, n)$
- Hash each member of $s \in S$ to one of n buckets, and set that bit to 1, i.e., $B[h(s)] = 1$
- Hash each element a of the stream and output only those that hash to bit that was set to 1
 - Output a if $B[h(a)] == 1$

First Cut Solution (2)



- **Creates false positives but no false negatives**
 - If the item is in S we surely output it, if not we may still output it

First Cut Solution (3)

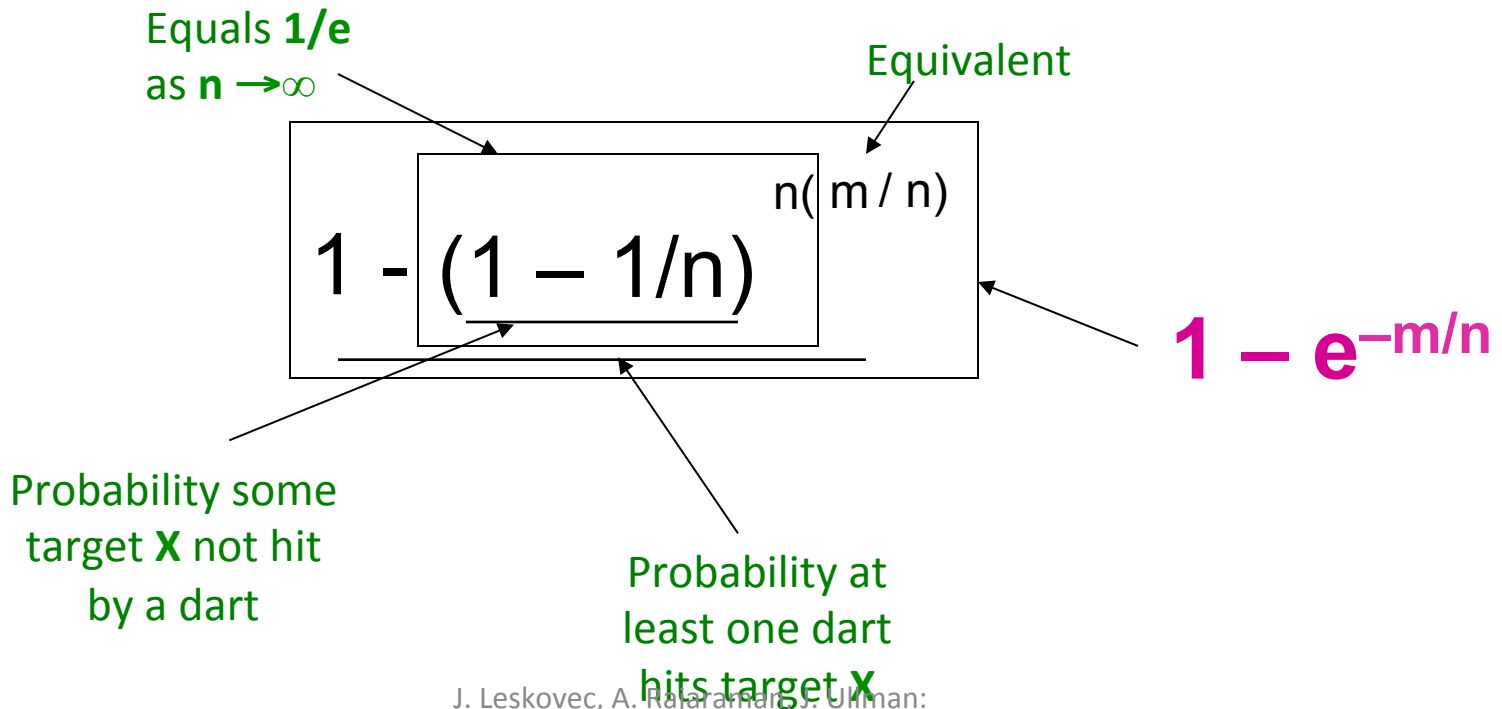
- $|S| = 1$ billion email addresses
 $|B| = 1\text{GB} = 8$ billion bits
- If the email address is in S , then it surely hashes to a bucket that has the bit set to **1**, so it always gets through (*no false negatives*)
- Approximately $1/8$ of the bits are set to **1**, so about $1/8^{\text{th}}$ of the addresses not in S get through to the output (*false positives*)
 - Actually, less than $1/8^{\text{th}}$, because more than one address might hash to the same bit

Analysis: Throwing Darts (1)

- **More accurate analysis for the number of false positives**
- **Consider:** If we throw m darts into n equally likely targets, **what is the probability that a target gets at least one dart?**
- **In our case:**
 - **Targets** = bits or buckets
 - **Darts** = hash values of items

Analysis: Throwing Darts (2)

- We have m darts, n targets
- **What is the probability that a target gets at least one dart?**



Analysis: Throwing Darts (3)

- **Fraction of 1s in the array \mathbf{B} =**
= probability of false positive = $1 - e^{-m/n}$
- **Example: 10^9 darts, $8 \cdot 10^9$ targets**
 - Fraction of **1s** in **\mathbf{B}** = **$1 - e^{-1/8} = 0.1175$**
 - Compare with our earlier estimate: **$1/8 = 0.125$**

Bloom Filter

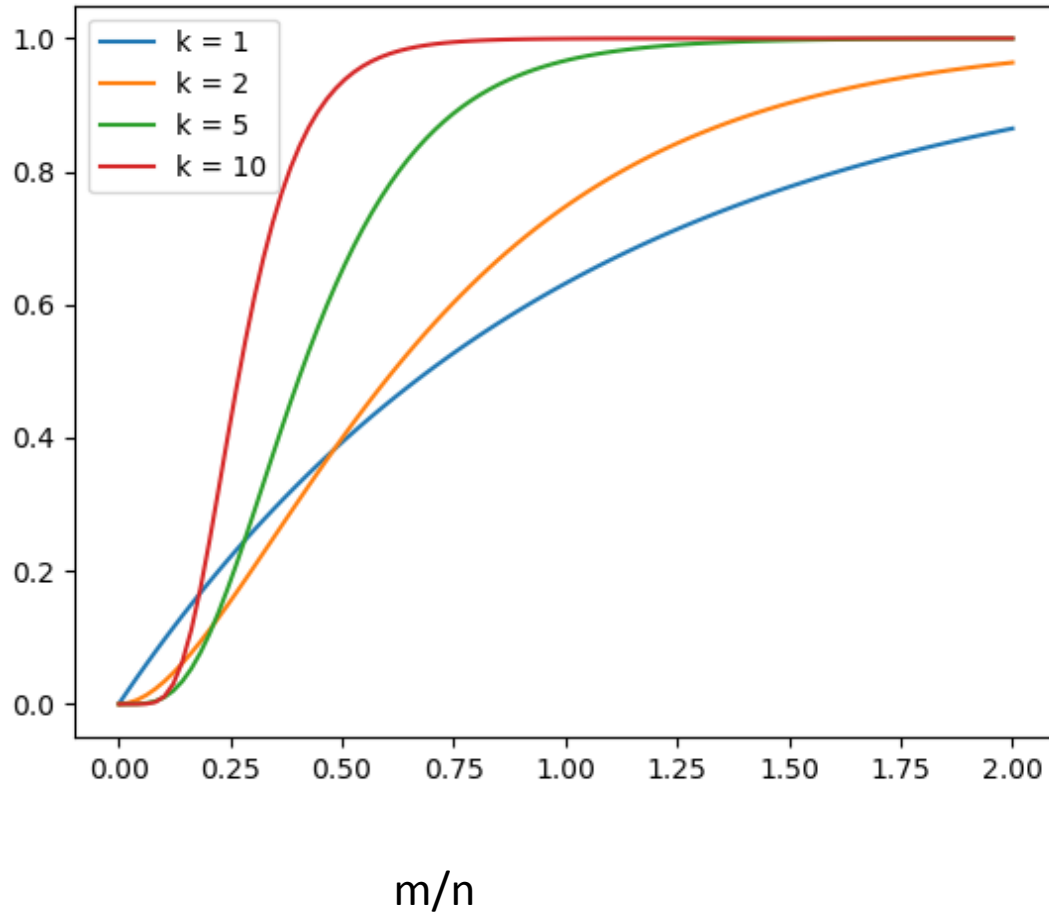
- Consider: $|\mathbf{S}| = m, |\mathbf{B}| = n$
- Use k independent hash functions h_1, \dots, h_k
- **Initialization:**
 - Set \mathbf{B} to all 0 s
 - Hash each element $s \in \mathbf{S}$ using each hash function h_i , set $\mathbf{B}[h_i(s)] = 1$ (for each $i = 1, \dots, k$) (note: we have a single array B!)
- **Run-time:**
 - When a stream element with key \mathbf{x} arrives
 - If $\mathbf{B}[h_i(\mathbf{x})] = 1$ for all $i = 1, \dots, k$ then declare that \mathbf{x} is in \mathbf{S}
 - That is, \mathbf{x} hashes to a bucket set to 1 for every hash function $h_i(\mathbf{x})$
 - Otherwise discard the element \mathbf{x}

Bloom Filter -- Analysis

- **What fraction of the bit vector B are 1s?**
 - Throwing $k \cdot m$ darts at n targets
 - So fraction of 1s is $(1 - e^{-km/n})$
- But we have k independent hash functions and we only let the element x through **if all k** hash element x to a bucket of value **1**
- So, false **positive probability** = $(1 - e^{-km/n})^k$

Bloom filter analysis

Prob. Of
false positive



Bloom Filter – Analysis (2)

- $m = 1$ billion, $n = 8$ billion

- $k = 1$: $(1 - e^{-1/8}) = 0.1175$

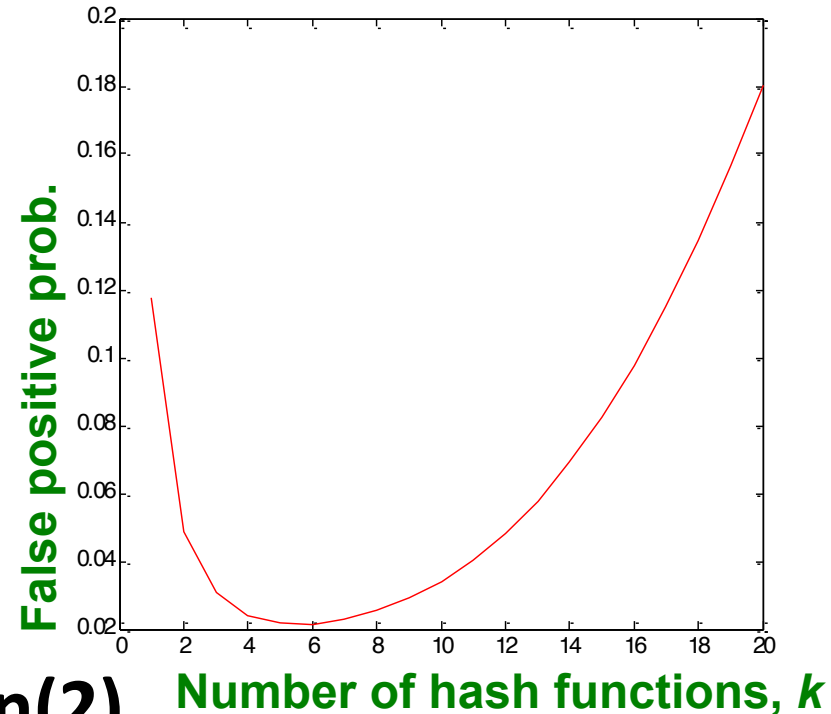
- $k = 2$: $(1 - e^{-1/4})^2 = 0.0493$

- What happens as we keep increasing k ?

- “Optimal” value of k : $n/m \ln(2)$

- In our case: Optimal $k = 8 \ln(2) = 5.54 \approx 6$

- Error at $k = 6$: $(1 - e^{-1/6})^2 = 0.0235$



Bloom Filter: Wrap-up

- **Bloom filters guarantee no false negatives, and use limited memory**
 - Great for pre-processing before more expensive checks
- **Suitable for hardware implementation**
 - Hash function computations can be parallelized
- **Is it better to have 1 big B or k small Bs?**
 - It is the same: $(1 - e^{-km/n})^k$ vs. $(1 - e^{-m/(n/k)})^k$
 - But keeping 1 big B is simpler

(2) Counting Distinct Elements

Counting Distinct Elements

- **Problem:**

- Data stream consists of a universe of elements chosen from a set of size N
- Maintain a count of the number of distinct elements seen so far

- **Obvious approach:**

Maintain the set of elements seen so far

- That is, keep a hash table of all the distinct elements seen so far

Applications

- **How many different words are found among the Web pages being crawled at a site?**
 - Unusually low or high numbers could indicate artificial pages (spam?)
- **How many different Web pages does each customer request in a week?**
- **How many distinct products have we sold in the last week?**

Using Small Storage

- **Real problem: What if we do not have space to maintain the set of elements seen so far?**
- **Estimate the count in an unbiased way**
- **Accept that the count may have a little error, but limit the probability that the error is large**

Flajolet-Martin Approach

- Pick a hash function h that maps each of the N elements to at least $\log_2 N$ bits
- For each stream element a , let $r(a)$ be the number of trailing 0s in $h(a)$
 - $r(a)$ = position of first 1 counting from the right
 - E.g., say $h(a) = 12$, then 12 is 1100 in binary, so $r(a) = 2$
- Record $R = \text{the maximum } r(a) \text{ seen}$
 - $R = \max_a r(a)$, over all the items a seen so far
- Estimated number of distinct elements = 2^R

Why It Works: Intuition

- Very very rough and heuristic intuition why Flajolet-Martin works:
 - $h(a)$ hashes a with **equal prob.** to any of N values
 - Then $h(a)$ is a sequence of $\log_2 N$ bits, where 2^{-r} fraction of all a s have a tail of r zeros
 - About 50% of a s hash to *****0**
 - About 25% of a s hash to ****00**
 - So, if we saw the longest tail of $r=2$ (i.e., item hash ending ***100**) then we have probably seen **about 4** distinct items so far
 - **So, it takes to hash about 2^r items before we see one with zero-suffix of length r**

Why It Works: More formally

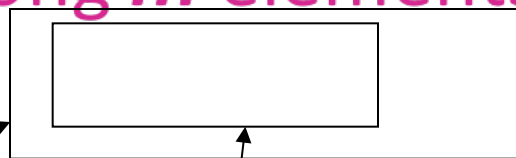
- Now we show why Flajolet-Martin works
- Formally, we will show that **probability of finding a tail of r zeros:**
 - Goes to **1** if $m \gg 2^{\uparrow r}$
 - Goes to **0** if $m \ll 2^{\uparrow r}$where m is the number of distinct elements seen so far in the stream
- Thus, 2^R will almost always be around m !

Why It Works: More formally

- What is the probability that a given $h(a)$ ends in at least r zeros is 2^{-r}
 - $h(a)$ hashes elements uniformly at random
 - Probability that a random number ends in at least r zeros is 2^{-r}
- Then, the probability of **NOT** seeing a tail of length r among m elements:

$$(1 - 2^{-r})^m$$

Prob. all end in fewer than r zeros.



Prob. that given $h(a)$ ends in fewer than r zeros

Why It Works: More formally

- **Note:** $(1 - 2^{-r})^m = (1 - 2^{-r})^{2^r (m2^{-r})} \approx e^{-m2^{-r}}$
- **Prob. of NOT finding a tail of length r is:**
 - If $m \ll 2^r$, then prob. tends to **1**
 - $(1 - 2^{-r})^m \approx e^{-m2^{-r}} = 1$ as $m/2^r \rightarrow 0$
 - So, the probability of finding a tail of length r tends to **0**
 - If $m \gg 2^r$, then prob. tends to **0**
 - $(1 - 2^{-r})^m \approx e^{-m2^{-r}} = 0$ as $m/2^r \rightarrow \infty$
 - So, the probability of finding a tail of length r tends to **1**
- **Thus, 2^R will almost always be around m !**

Why It Doesn't Work

- **$E[2^R]$ is actually infinite**
 - Probability halves when $R \rightarrow R+1$, but value doubles
- **Workaround involves using many hash functions h_i and getting many samples of R_i**
- **How are samples R_i combined?**
 - **Average?** What if one very large value $2^{\uparrow R} \downarrow i$?
 - **Median?** All estimates are a power of 2
 - **Solution:**
 - Partition your samples into small groups
 - Take the median of groups
 - Then take the average of the medians

(3) Computing Moments

Generalization: Moments

- Suppose a stream has elements chosen from a set A of N values
- Let m_i be the number of times value i occurs in the stream
- The k^{th} *moment* is

$$\sum_{i \in A} (m_i)^k$$

Special Cases

$$\sum_{i \in A} (m_i)^k$$

- **0th moment** = number of distinct elements
 - The problem just considered
- **1st moment** = count of the numbers of elements = length of the stream
 - Easy to compute
- **2nd moment** = *surprise number S* = a measure of how uneven the distribution is

Example: Surprise Number

- **Stream of length 100**
- **11 distinct values**
- Item counts: **10, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9**
Surprise $S = 910$
- Item counts: **90, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1**
Surprise $S = 8,110$

AMS method

- AMS method works for all moments
- Gives an unbiased estimate.
- We will just concentrate on the 2nd moment S .
- We pick and keep track of many variables X :
 - For each variable X , store $X.el$ and $X.val$
 - $X.el$ corresponds to the item l
 - $X.val$ corresponds to the count of item l
 - Note this requires a count in main memory, so number of X s is limited
- Our goal is to compute $S = \sum_i m_i^2$

One random variable (X)

- How to set $X.val$ and $X.el$?
 - Assume stream has length n (we relax this later)
 - Pick some random time t ($t < n$) to start, so that any time is equally likely
 - Let at time t the stream have item i . We set $X.el = i$
 - Then we maintain count c ($X.val = c$) of the number of i s in the stream starting from the chosen time t

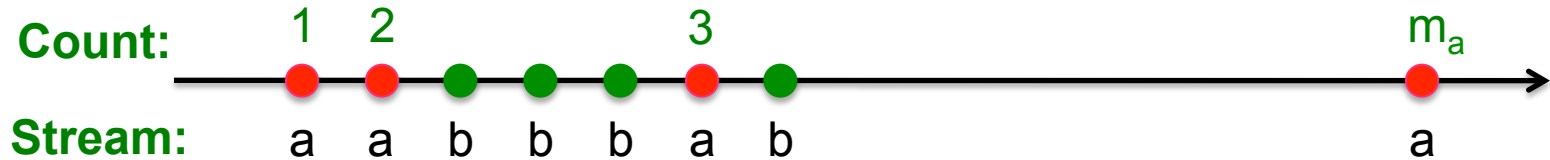
- Then the estimate of the 2nd moment ($\sum_i m_i^2$) is:

$$S = f(X) = n(2c - 1)$$

- Note, we will keep track of multiple X s, (X_1, X_2, \dots, X_k) and our final estimate will be:

$$S = 1/k \sum_j f(X_j)$$

Expectation Analysis



- 2nd moment is $S = \sum_i m_i^2$
- C_t - number of times item at time t appears from time t onwards ($c_1=m_a$, $c_2=m_a-1$, $c_3=m_b$)
- $E[f(X)] = 1/n \sum_{t=1}^n n (2c_t - 1)$
 $= 1/n \sum_i n (1 + 3 + 5 + \dots + 2 m_i - 1)$

m_i ... total count of item i in the stream (we are assuming stream has length n)

Group times by the value seen

Time t when the last i is seen ($c_t=1$)

Time t when the penultimate i is seen ($c_t=2$)

Time t when the first i is seen ($c_t=m_i$)

Higher-Order Moments

- For estimating k th moment we essentially use the same algorithm but change the estimate:
 - For $k=2$ we used $n(2 \cdot c - 1)$
For $k=3$ we use: $n(3 \cdot c^2 - 3c + 1)$ (where $c=X.val$)
- Why?
 - For $k=2$: Remember we had $(1+3+5+\dots+(2m_i-1))$ and we showed terms $2c-1$ (for $c=1,\dots,m$) sum to m^2
 - $2c - 1 = c^2 - (c-1)^2$
 - For $k=3$: $c^3 - (c-1)^3 = 3c^2 - 3c + 1$
- Generally: Estimate = $n(c^k - (c-1)^k)$

Combining Samples

- **In practice:**
 - Compute $f(\mathbf{X}) = \mathbf{n}(2\mathbf{c} - \mathbf{1})$ for as many variables \mathbf{X} as you can fit in memory
 - Average them in groups
 - Take median of averages
- **Problem: Streams never end**
 - We assumed there was a number \mathbf{n} , the number of positions in the stream
 - But real streams go on forever, so \mathbf{n} is a variable – the number of inputs seen so far

Streams Never End: Fixups

- **(1)** The variables X have n as a factor – keep n separately; just hold the count in X
- **(2)** Suppose we can only store k counts. We must throw some X s out as time goes on:
 - **Objective:** Each starting time t is selected with probability k/n
 - **Solution: (fixed-size sampling!)**
 - Choose the first k times for k variables
 - When the n^{th} element arrives ($n > k$), choose it with probability k/n
 - If you choose it, throw one of the previously stored variables X out, with equal probability

AMS ALGORITHM USING SKETCHES

Generalization of AMS Algorithm

- Stream of pair (i,c) , $i \in \{1, \dots, U\}$ and c is positive integer.
- $x[i] = x[i] + c$ for each update
- Join size: $\mathbf{x} \cdot \mathbf{y} = \sum_{i=1}^U (x[i] y[i])$
- Pth Moment: $F_p(\mathbf{x}) = \sum_{i=1}^U x[i]^p$

$$\|\mathbf{x} - \mathbf{y}\|_2 = \sqrt{F_2(\mathbf{x} - \mathbf{y})}.$$

- $h : \{1, \dots, U\} \rightarrow \{+1, -1\}$

Generalization of AMS Algorithm

UPDATE(i, c, z)

Input: item i , count c , sketch z

- 1: **for** $j = 1$ to w **do**
- 2: **for** $k = 1$ to d **do**
- 3: $z[j][k] += h_{j,k}(i) * c$

ESTIMATE $F_2(z)$

Input: sketch z

- 1: **Return** ESTIMATEJS(z, z)

ESTIMATEJS(x, y)

Input: sketch x , sketch y

Output: estimate of $x \cdot y$

- 1: **for** $j = 1$ to w **do**
- 2: $avg[j] = 0;$
- 3: **for** $k = 1$ to d **do**
- 4: $avg[j] += x[j][k] * y[j][k] / w;$
- 5: **Return**(median(avg))

Fig. 1 AMS algorithm for estimating join and self-join size

Generalization of AMS Algorithm

Lemma 1 $E(Z^2) = F_2(\mathbf{x})$

Proof

$$\begin{aligned} E(Z^2) &= E\left(\left(\sum_{i=1}^U h(i)\mathbf{x}[i]\right)^2\right) \\ &= E\left(\sum_{i=1}^U h(i)^2\mathbf{x}[i]^2\right) + E\sum_{1 \leq i < j \leq U} 2h(i)h(j)\mathbf{x}[i]\mathbf{x}[j] \\ &= \sum_{i=1}^U \mathbf{x}[i]^2 + 0 = F_2(\mathbf{x}). \end{aligned}$$

Generalization of AMS Algorithm

- $\text{Var}(Z^2) \leq 2F_2(\mathbf{x})^2$

$$\begin{aligned}\text{Var}(Z^2) &= \mathbb{E}(Z^4) - \mathbb{E}(Z^2)^2 \\ &= \mathbb{E}\left(\left(\sum_{i=1}^U h(i)\mathbf{x}[i]\right)^4\right) - \left(\sum_{i=1}^U \mathbf{x}[i]^2\right)^2\end{aligned}$$

Generalization of AMS Algorithm

$$\begin{aligned}
 &= \mathbb{E} \left(\left(\sum_{i=1}^U h(i)^4 \mathbf{x}[i]^4 + \sum_{1 \leq i < j \leq U} 6h(i)^2 h(j)^2 \mathbf{x}[i]^2 \mathbf{x}[j]^2 \right. \right. \\
 &\quad + \sum_{i, i \neq j \neq k} 12h(i)^2 h(j) h(k) \mathbf{x}[i]^2 \mathbf{x}[j] \mathbf{x}[k] \\
 &\quad + \sum_{1 \leq i \neq j \leq U} 4h^3(i) h(j) \mathbf{x}[i]^3 \mathbf{x}[j] \\
 &\quad \left. \left. + \sum_{1 \leq i < j < k < l \leq U} 12h(i) h(j) h(k) h(l) \mathbf{x}[i] \mathbf{x}[j] \mathbf{x}[k] \mathbf{x}[l] \right) \right) \\
 &\quad - \left(\sum_{i=1}^U \mathbf{x}[i]^4 + \sum_{1 \leq i < j \leq U} 2\mathbf{x}[i]^2 \mathbf{x}[j]^2 \right) \\
 &\quad \dots
 \end{aligned}$$

Generalization of AMS Algorithm

$$\begin{aligned} &= \sum_{i=1}^U \mathbf{x}[i]^4 + \sum_{1 \leq i < j \leq U} 6\mathbf{x}[i]^2 \mathbf{x}[j]^2 \\ &\quad - \left(\sum_{i=1}^U \mathbf{x}[i]^4 + \sum_{1 \leq i < j \leq U} 2\mathbf{x}[i]^2 \mathbf{x}[j]^2 \right) \\ &= 4 \sum_{1 \leq i < j \leq U} \mathbf{x}[i]^2 \mathbf{x}[j]^2 \leq 2F_2^2. \end{aligned}$$

Generalization of AMS Algorithm

Fact 1 (Variance Reduction) *Let X_i be independent and identically distributed random variables. Then*

$$\text{Var}\left(\sum_{i=1}^w \frac{X_i}{w}\right) = \frac{1}{w} \text{Var}(X_1).$$

Fact 2 (The Chebyshev Inequality) *Given a random variable X ,*

$$\Pr[|X - \mathbb{E}(X)| \geq k] \leq \frac{\text{Var}(X)}{k^2}.$$

Generalization of AMS Algorithm

Theorem 1 *An (ϵ, δ) -approximation of F_2 , the self-join size, can be computed in space $O(\frac{1}{\epsilon^2} \log 1/\delta)$ machine words in the streaming model. Each update takes time $O(\frac{1}{\epsilon^2} \log 1/\delta)$.*

Proof Applying the Chebyshev inequality to the average of $w = \frac{16}{\epsilon^2}$ copies of the estimate Z generates a new estimate Y such that

$$\Pr[|Y - F_2| \leq \epsilon F_2] \leq \frac{\text{Var}(Y)}{\epsilon^2 F_2^2} = \frac{\text{Var}(Z)}{c\epsilon^2 F_2^2} = \frac{2F_2^2}{(16/\epsilon^2)\epsilon^2 F_2^2} = \frac{1}{8}.$$

Generalization of AMS Algorithm

Fact 3 (Application of Chernoff Bounds) *Let R be a range of values $R = [R_{\min}..R_{\max}]$, and let Y_i be $d = 4 \log 1/\delta$ independent and identically distributed random variable such that $\Pr[Y_i \notin R] \leq \frac{1}{8}$. Then*

$$\Pr[(\text{median}_{i=1}^d Y_i) \notin R] \leq \delta,$$

that is, if there is constant probability that each Y_i falls within the desired range R , then taking the median of $O(\log 1/\delta)$ copies of Y_i reduces the failure probability to δ .

Hence, applying the Chernoff bound result from Fact 3 to the median of $4 \log 1/\delta$ copies of the average Y gives the probability of the results being outside the range of ϵF_2 from F_2 as δ . The space required is that to maintain $O(\frac{1}{\epsilon^2} \log 1/\delta)$ copies of the original estimate. Each of these requires a counter and a 4-wise independent hash function, both of which can be represented with a constant number of machine words under the standard RAM model. □

Join size estimation

Lemma 3 *Let Z_x be an entry of a sketch computed for the vector \mathbf{x} , and let Z_y be an entry of a sketch computer for \mathbf{y} using the same hash function. The estimate is correct in expectation, i.e., $\mathbb{E}(Z_x * Z_y) = \mathbf{x} \cdot \mathbf{y}$.*

Proof

$$\begin{aligned}\mathbb{E}(Z_x * Z_y) &= \mathbb{E}\left(\sum_{i=1}^U h(i)^2 \mathbf{x}[i] \mathbf{y}[i] + \sum_{1 \leq i \neq j \leq U} h(i) h(j) \mathbf{x}[i] \mathbf{y}[j]\right) \\ &= \sum_{i=1}^U \mathbf{x}[i] \mathbf{y}[i] + 0 = \mathbf{x} \cdot \mathbf{y}.\end{aligned}$$

□

Join size estimation

$$\text{Var}(Z_x * Z_y) \leq F_2(\mathbf{x})F_2(\mathbf{y}).$$

Theorem 2 *Using space $O(\frac{1}{\epsilon^2} \log 1/\delta)$ space we can output an estimate of $\mathbf{x} \cdot \mathbf{y}$ so that*

$$\Pr[|(\mathbf{x} \cdot \mathbf{y}) - est| \leq \epsilon \sqrt{F_2(\mathbf{x})F_2(\mathbf{y})}] \geq 1 - \delta.$$