# L12: Streaming : Count-Min Sketch and Others

Jeff M. Phillips

February 21, 2018

- Count-Min Stretch
  └▸ proof
- Count Stretch
- Fregoent Itruset (Apriori Alg)
- Bloom Filters

# Streaming Model

$a_i \in [m]$

Input $A = \langle a_1, a_2, \ldots a_i, \ldots a_n \rangle$

$$\underbrace{\qquad}_{A_i}$$

$A_i := \{a_1, a_2, \ldots a_i\}$



Memory

m, n very large
$\rightarrow$ use $O(\log n + \log m)$
space

Frequency $f_j = |\{a_i \in A \mid a_i = j\}|$
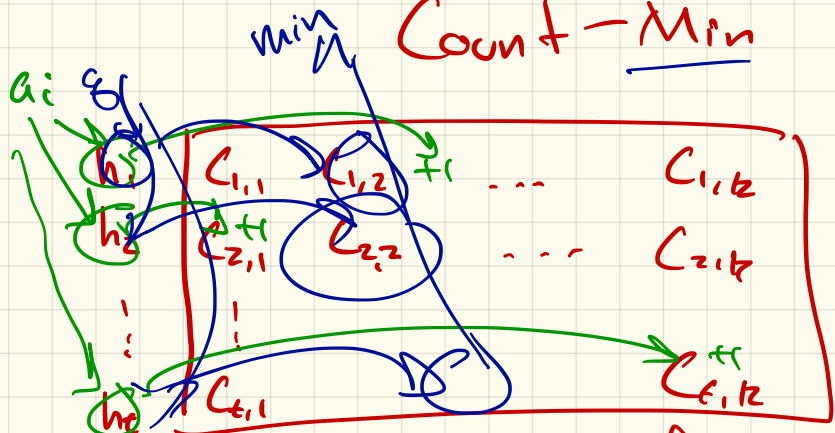
MG: $f_j - \varepsilon n \leq \hat{f}_j \leq f_j$

CM: $f_j \leq \hat{f}_j \leq f_j + \varepsilon n$

$\underbrace{\qquad}_{\text{hold w.p. } 1-\delta}$

$\leftarrow$ also handle
substractions
"turnstile"

# Count-Min Sketch



$t \cdot k$ counters

$k = 2/\varepsilon$

$t = \log(1/\delta)$

$+$ hash fxns

$h_j : [m] \to [k]$

$\text{space}(C_{ij}) = O(\log n)$

$\text{space}(h_j) = O(\log m)$

**Initialize** $C_{ij} = 0 \quad \forall i, j$

for $a_i \in A$
  for $j = 1$ to $t$
    $C_{j, h_j(a_i)} = C_{j, h_j(a_i)} + 1$

Query $\hat{f}_g$ ? $g \in [m]$

$\hat{f}_g = \min_{j' \in [t]} C_{j, h_j(g)}$

$$\underline{Clear} \qquad f_g \leq \hat{f}_g \qquad \leftarrow only \ overcounts$$

$$\hat{f}_g \overset{!}{\leq} f_g + w \qquad\qquad p \in [m]$$

R.V. $\quad Y_{p,j} = \begin{cases} f_p & w.p. \ 1/k \\ 0 & otherwise \end{cases}$ $\Big)$ $j$th row

Prob $c_{j, h_j(g)}$

the overcount from

$p \in [m]$
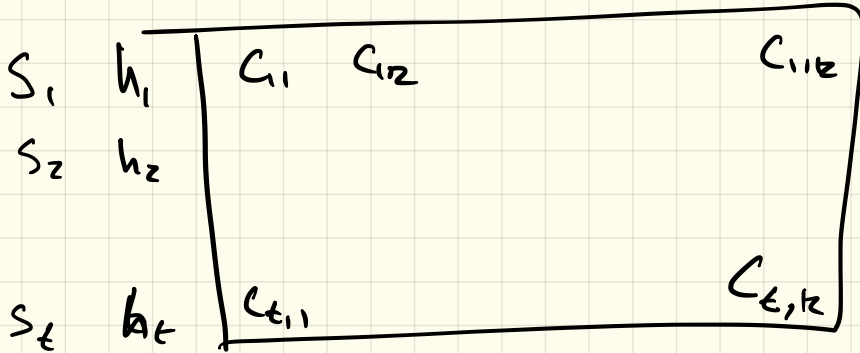
$$E[Y_{p,j}] = f_p / k$$

R.V. $\quad X_j = \displaystyle\sum_{\substack{p \in [m] \\ \neq g}}' Y_{p,j} \quad \leftarrow$ in $j$th row, total

over count on

$c_{j, h_j(g)}$

$$E[X_j] = E\left(\sum_{p \neq g}' Y_{p,j}\right) = \sum_{p \neq g}' f_p/k \leq \frac{n}{k} = \frac{\varepsilon n}{2}$$

$$Pr[X > \alpha] = \frac{E[\cancel{X}]}{\alpha} = \frac{1}{2} \quad \Big| \quad Prob \ all \ t \ rows \overset{> \varepsilon n \ error}{(1/2)^t}$$

$$\alpha = E[X] \cdot 2$$

# Count Sketch

$$\begin{array}{cc}
S_1 & h_1 \\
S_2 & h_2 \\
\\
S_t & h_t
\end{array}$$

$$\begin{array}{ccccc}
C_1 & C_{12} & & & C_{1,k} \\
\\
\\
C_{t,1} & & & & C_{t,k}
\end{array}$$

$$\boxed{k = \frac{1}{\varepsilon^2}}$$

$$t = \log\left(\frac{2}{\delta}\right)$$

$$h_j : [m] \to [k] \quad (\text{random}) \qquad \overset{(\text{random})}{S_j : [m] \to \{-1, +1\}}$$

$$E[C_{ij}] = 0$$

for $a_i \in A$
    for $j \in [t]$
        $$C_{j, h_j(a_i)} = C_{j, h_j(a_i)} + S_j(a_i) \cdot 1$$

$$\left| f_g - \hat{f}_g \right| \leq \varepsilon F_2$$

$$F_2 = \sqrt{\sum_{p \in [m]} f_p^2}$$

# Bloom Filter

Data Structure $S$ for sets.

Streams for $a_i \in A$

$$\text{put } a_i \xrightarrow[\text{into}]{} S$$

---

Query is $q \in [m]$ in $S$?

- if $q \in S \implies$ <u>always</u> return <u>true</u>
- if $q \notin S \implies$ usually return <u>false</u>

Init $B[q] = 0$ $\forall q$

for $a_i \in A$

for $j = 1$ to $k$
  Set $B[h_j(a_i)] = 1$

$k$ hash funs $h_1, h_2, \ldots h_k$

$1$ arrays of bits of $m$ bits $B[\ ]$

$$k \approx \frac{m}{n} \ln(2)$$

# A-Priori Algorithm (Frequent Itemsets)

Input: $A = \{a_1, a_2 \ldots a_n\}$

$$a_i = \{x_1, x_7, x_{14}\} \subset [m]$$

Market Basket Analysis

$\hookrightarrow$ beer + diapers

$\varepsilon = 0.05$

Find all tuples $\{x_1, x_2, x_2\}$ w/ cooccur in at least $\boxed{\varepsilon n}$ baskets

If $\{x_1, x_2, x_7\}$ cooccur in 5% then each of $x_1, x_7,$ and $x_2$ most each occur in 5%

# Frequent Itemsets : Apriori

Find tuples in at least 1/3 sets

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 3 | 3 | 5 | 4 | 3 | 3 | 8 | 4 | 2 | 4 |

| 2,3 | 2,6 | 2,7 | 2,9 | 3,6 | 3,7 | 3,9 |
|-----|-----|-----|-----|-----|-----|-----|
| 1 | 3 | 1 | 2 | 3 | 0 | 1 |

| 6,7 | 1,9 | 7,9 |
|-----|-----|-----|
| 3 | 4 | 2 |

3, 6, 9

$T_1 = \{1, 2, 3, 4, 5\}$

$T_2 = \{2, 6, 7, 9\}$

$T_3 = \{1, 3, 5, 6\}$

$T_4 = \{2, 6, 9\}$

$T_5 = \{7, 8\}$

$T_6 = \{1, 2, 6\}$

$T_7 = \{0, 3, 5, 6\}$

$T_8 = \{0, 2, 4\}$

$T_9 = \{2, 4\}$

$T_{10} = \{6, 7, 9\}$

$T_{11} = \{3, 6, 9\}$

$T_{12} = \{6, 7, 8\}$