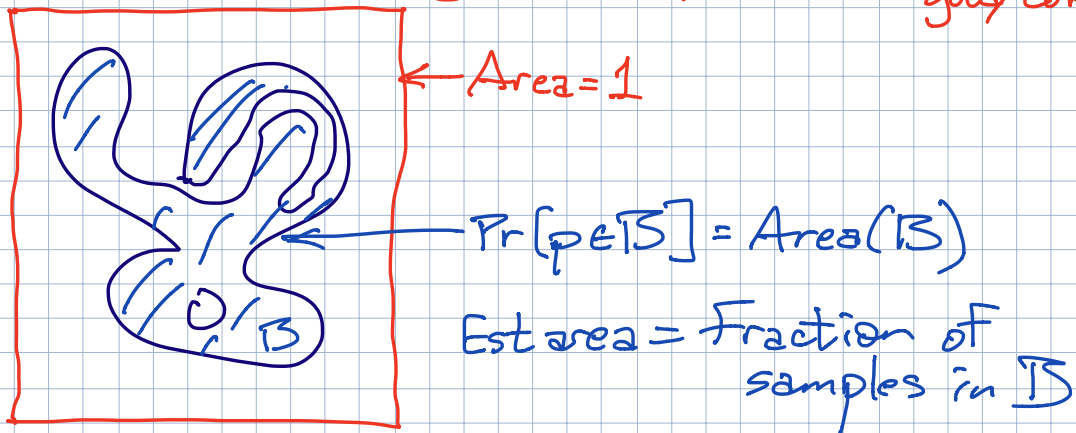


HW4 due Tue Mar 1

Randomized algo's with prob. of error

Las Vegas algo: Always correct, probably fast

Monte Carlo algo: Always fast, probably good/correct



Monte Carlo: error rate δ confidence $1 - \delta$

$$\Pr[\text{wrong/bad}] < \delta$$

Set membership: Insert Query
with error δ

$x \in S \Rightarrow \text{YES}$

$x \notin S \Rightarrow \text{NO}$ with prob $1 - \delta$

one-sided error

YES w/ prob $< \delta$

FILTER

ZATOCODING '57

Bloom Filter '70

Simultaneously Selective Patterns

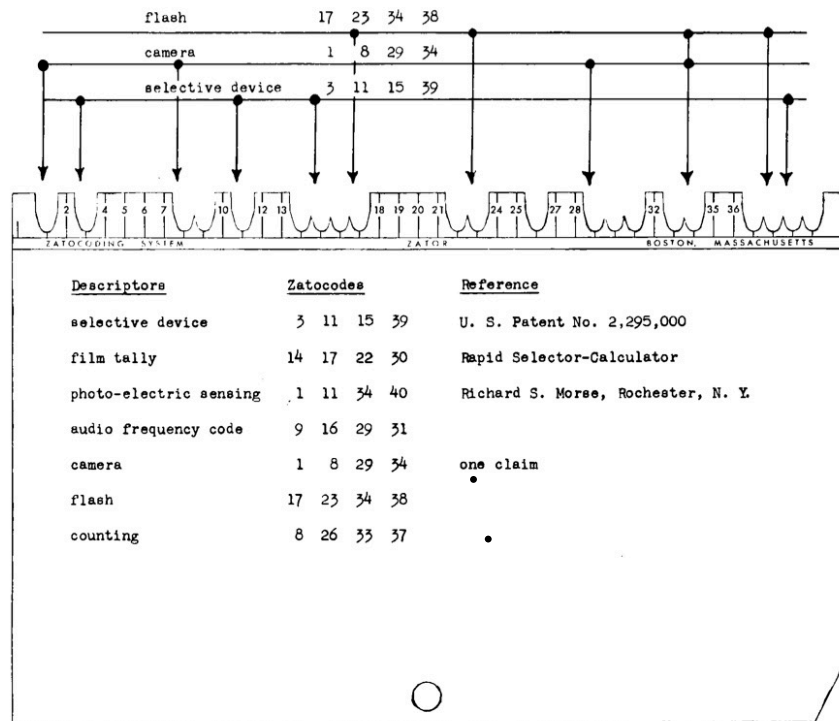


Fig. 1. ZATOCODING, illustrated with a 5 by 8 inch edge notched ZATOCARD for the tally. Note that the random ZATOCODE patterns in the field overlap and intermingle. Selection on the combination of three descriptors, "flash," "camera," and "selective device," is according to the inclusion of the pattern of arrows into the pattern of notches in the coding field. ZATOCARDS are sorted by the selector shown in Figure 2.

Bloom filter: bits $B[0..m-1]$

+ k hash functions $h_i: \mathcal{U} \rightarrow [m]$

ideal random

Insert(x):

for $i \leftarrow 1$ to k
 $B[h_i(x)] \leftarrow 1$

Member?(x):

for $i \leftarrow 1$ to k
if $B[h_i(x)] = 0$
FALSE
TRUE

$$\Pr[h_i(x) = j] = 1/m$$

$$\Pr(\text{Insert}(x) \text{ does not set } B[j] \leftarrow 1) = \left(1 - \frac{1}{m}\right)^k$$

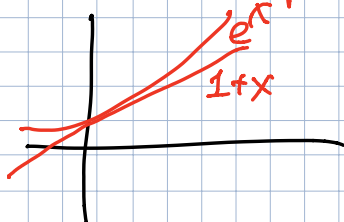
After n Inserts:

$$\Pr[B[j] = 0] = \left(1 - \frac{1}{m}\right)^{kn} < e^{-kn/m}$$

WMU \neq

$$e^x \geq 1+x$$

$$\Pr[\text{False positive}] = (1-p)^k \quad (\text{sort of})$$



$$\delta = (1 - e^{-kn/m})^k$$

$$\ln \delta = \ln(1-p)^k = k \ln(1-p) = -\frac{m}{n} \ln p \ln(1-p)$$

$$\text{max at } \boxed{p=1/2} \Rightarrow \boxed{k = \frac{m}{n} \cdot \ln 2}$$

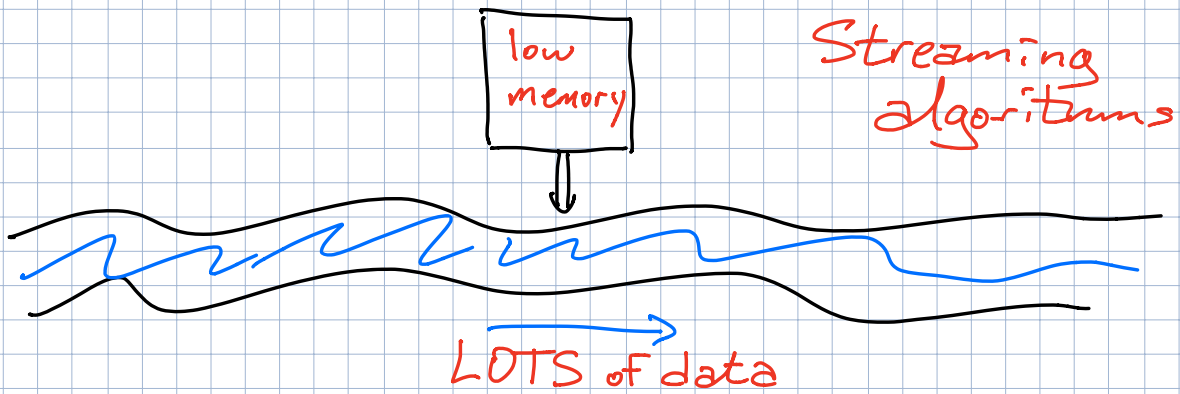
$$\delta = \left(\frac{1}{2}\right)^{\frac{m}{n} \ln 2} \approx (0.6185)^{m/n}$$

$$\boxed{m = \frac{\log(1/\delta)}{\ln 2} \cdot n} \Rightarrow \text{error prob } \delta$$
$$= O(n \cdot \log(1/\delta))$$

$$\delta = 1\% \Rightarrow 10n \text{ bits } \quad k = 7$$

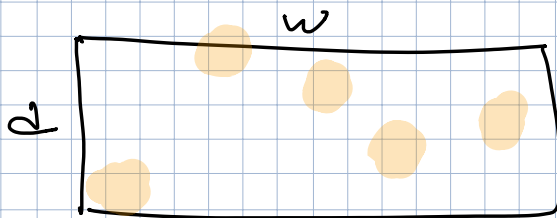
$$m = 32n \Rightarrow k = 22 \quad \delta \approx 2 \cdot 10^{-7}$$

Mittemacher
Broder]



"How many times have I seen item x ?"

Count Min Sketch $w \times d$ array of integers
width depth universal



d hash functions
 $h_i: \mathcal{U} \rightarrow [w]$

Process(x):

For $i \leftarrow 1$ to d
 $\text{Count}[i, h_i(x)]++$

Estimate(x):

$\min_i \text{Count}[i, h_i(x)]$

If we choose w and d correctly

$$\Pr[\text{Estimate}(x) > \text{freq}(x) + \epsilon \cdot N] \leq \delta$$

total length of stream so far

N

$$w = \left\lceil \frac{e}{\epsilon} \right\rceil \quad d = \lceil \ln(1/\delta) \rceil$$

$E[X_{i,x}] = \# \text{ collisions with } x \text{ in row } i$

$$= \sum_{y \neq x} \Pr[h_i(x) = h_i(y)] \cdot \text{Freq}(y)$$

$$\leq \frac{N}{w}$$

Markov's inequality.

$$\Pr[X_{i,x} > \epsilon N] \leq \frac{1}{w\epsilon}$$

$$\Pr[\text{Est} > \text{Freq}(x) + \epsilon N] \leq \left(\frac{1}{w\epsilon}\right)^d < \delta$$