

RANDOM SAMPLING:

$\rightarrow \mathbb{N}, \mathbb{R}, \text{strings} \dots$

- Given a sequence of items $S = (i_1, i_2, \dots, i_m)$ and a positive integer $m \leq M \Rightarrow$ select a subset of m items from S uniformly at random

1st Disk streaming can't set prob. in advance
 if n is known (or not) if n is known \Rightarrow if know prob. of extraction in advance

if n is known $\Rightarrow P(\text{sampling an item}) = \frac{1}{m}$

items occupy $\frac{m}{B}$ pages

- procedure $\text{Rand}(A, B)$: selects at random a number in range $[a, b]$

- want position of sampled items in sorted order: P2 CSES : m known or m unknown

- \Rightarrow 1) speeds up the extraction from S (Disk & Stream)
 2) Reduces working space

- extracts items efficiently via scan (no auxiliary array of pointers)

if $m > M \Rightarrow$ need a disk-based sorter

if $m \leq M \Rightarrow$ positions are integers in a fixed range

\Rightarrow radix sort or others

- ① m is known and 2 LEVEL MEMORY MODEL
 DISK MODEL and known sequence length

- input size m is known

- $S[1, m]$ stored in a file on disk and cannot be modified

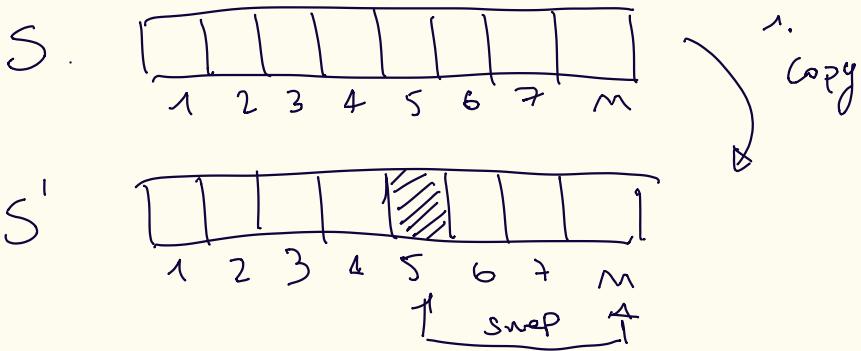
ALGORITHM

- initialize the auxiliary array $S'[1, m] = S[1, m];$
- for $s = 0, 1, \dots, m-1$ do
- $p = \text{Rand}(1, m-s);$
- select the item (pointed by) $S'[p];$
- swap $S'[p]$ with $S'[m-s]$
- end for

of pointers (wz: items may link variable length, os. strings)

// copy $S \rightarrow S'$

$s = 0$ ultimes elements



2. $s \sim$

of sampled items

3. $p = \text{rand}(1, m-s) \Rightarrow$ item to be sampled

$$p: \text{rand}(1, 8) \Rightarrow$$

$$= 5$$

u: pick $S'[s]$

5. swap $S'[s]$ with last item

$$S'[m-s] :=$$

$$S'[8-n]$$

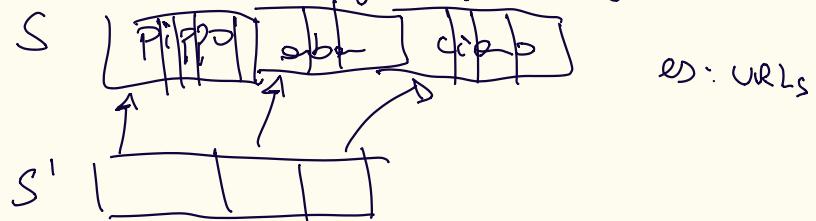
Repeat the process:

$s=1 \Rightarrow$ sample between 1 and $M-s = M-1$

GOOD

- # of RANDOM numbers is M • times good
- BAD
- ADDITIONAL space (I'm copying the array)
- lot of I/Os, I'm jumping

If m & S the arbitrary length strings



\Rightarrow In S' I can have pointers

$\Rightarrow S'$ is proportional to the # of objects, not total size of the set I wanna sample

$|S'| = n = \# \text{ items}$ (not to their length)

$$S = [31, 37, 53, 47, 97, 89, 11, 2, 3, 5] \quad m=4$$

$$S' = S;$$

for $s=0$ to $m-1$ do

$$p = \text{rand}(0, m-s)$$

$$a = S'[p]$$

swap $S'[p]$ with $S'[m-s-1]$

$$\begin{aligned} \text{1st RWN} & \quad \text{FOR } s=0 \quad \text{TO } 4-1=3 \\ & \quad p = \text{RAND}(0, 10-0) = 9 \quad \text{1st Run } = 5 \\ & \quad a = S'[9] \\ & \quad \text{SWAP } S'[9] \underset{s}{=} 5 \text{ with } S'[10-0-1] = S'[9] \underset{s}{=} 5 \end{aligned}$$

$$S = 1 \quad \text{TO } 3 \quad 9$$

$$p = \text{rand}(0, 10-1) = 4$$

$$a = S'[4] = 87$$

$$\text{SWAP } S'[4] \text{ with } S'[10-1-1] = S'[8] = 3$$

— — —

$$S = 2 \quad 10 \quad 3$$

$p \neq \text{rand} \dots$



Solution 2 (still n is known, avoiding jumps on disks)

- access items via scan-based access

\Rightarrow assume I can grab the index of the sampled items
(e.g. 20, 7, 5, 12)

\Rightarrow sort the indexes

\Rightarrow bring them sorted

\Rightarrow access from left to right \Rightarrow still $O(m)$ I/Os

but better w.r.t. main memory access, head of disk is faster than go back and forth

\Rightarrow find set of indexes for we sample?

list, hash table ...

Algorithm (using a dictionary DS)

- ~~AVoids items swapping via an auxiliary array~~

Algorithm 3.2 Dictionary of sampled positions

```
1: Initialize the dictionary  $D = \emptyset$ 
2: while ( $|D| < m$ ) do
3:    $p = \text{Rand}(1, n)$ ;
4:   if  $p \notin D$  insert it;
5: end while
```

assuming $m \leq n$

$\Rightarrow D$ stays

in memory

I/O free

I'm not touching the disk

SORT D

extract D 's items from S (disk) (left to right)

cost:

1) unsuccessful (extract a position already in D)

2) successful (.. $\in D$)

• CAN ASSUME
 $m \leq m/2$

How many time unsuccessful

$$P(p \in D) = \frac{\text{positive events}}{\text{total # of events}} = \frac{|D|}{m} < \frac{m}{m} \leq \frac{m/2}{m} = \frac{1}{2}$$

UNSUCCESSFUL \nearrow
Space from which
I sample p

- we can assume that $m < \frac{1}{2}n$ w.r.t. if $m > \frac{1}{2}n$ can turn
upside down the problem: in D I put the position I don't want to sample
(negative stuff) (complement)

• time complexity $O(m)$ time and $O(m)$ additional space

FACT 3.1 Algorithm 3.2 based on hashing with chaining requires $O(m)$ average time and takes $O(m)$ additional space to select uniformly at random m positions in $[1, n]$. The average depends both on the use of hashing and the cost of re-sampling. An additional sorting-cost is needed if we wish to extract the sampled items of S in a streaming-like fashion. In this case the overall sampling process takes $O(\min\{m, n/B\})$ I/Os.

if use a binary-search tree $O(m \log m)$ time

SOLUTION ③ avoid D (related to birthday paradox)

1. extract m numbers in $[1, n]$

if they are all distinct \rightarrow sort them

else REPEAT

checked by sorting

STREAMING MODEL, m known

Algorithm 3.4 Scanning and selecting

```

1:  $s = 0$ ;
2: for ( $j = 1; (j \leq n) \&& (s < m); j++$ ) do
3:    $p = \text{Rand}(0, 1)$ ;
4:   if ( $p < \frac{m-s}{n-j+1}$ ) then
5:     select  $S[j]$ ;
6:      $s++$ ;
7:   end if
8: end for

```

$s = 0$ m
 $j = 0; 0 \leq j \leq 2$

$$p = 0.4$$

$$\text{if } 0.4 \leq \frac{2-0}{2-0+1} = \frac{2}{3} \text{ then}$$

select $S[0]$

$$S = \begin{pmatrix} 1, & 12, & 14 \\ 0 & 1 & 2 \end{pmatrix}$$

want 2 samples: n_m

$s = 0$

for ($j = 0; j \leq m \&\& s < m; j++$)

$p = \text{rand}(0, 1)$;

if ($p < \frac{m-s}{m-j+1}$)

select $S[j]$

$s++$

end if
end for

$s = 1$ $m = 2$

$$j = 1; 1 \leq 3; 1 \leq 2$$

$$p = 0.5 < \frac{2-1}{3-1+1} = \frac{1}{3} \text{ if } 0.5 \text{ no } 2 \text{ not picked}$$

SOLUTION (h)

• m is known

• I'm considering STREAMING MODEL

\Rightarrow means S is looked from left to right

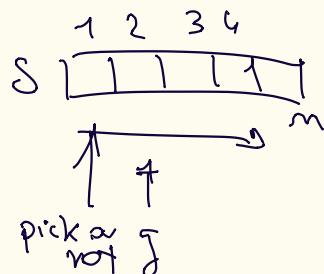
(no dictionary of position, may have not much space in memory, or search engines)

\Rightarrow look item i (decide pick or not)

go forward

2 (pick or not)

\Rightarrow NETWORK SITUATION: buffer only items I really pick

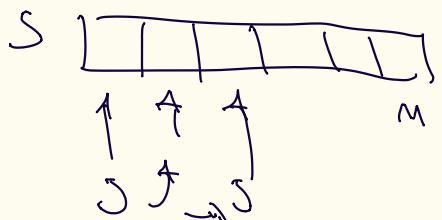


Algorithm 3.4 Scanning and selecting

```

1:  $s = 0$ ; # of sampled items
2: for ( $j = 1$ ; ( $j \leq n$ )  $\&\&$  ( $s < m$ );  $j++$ ) do
3:    $p = \text{Rand}(0, 1)$ ; Real
4:   if ( $p < \frac{m-s}{n-j+1}$ ) then items are less frequent
5:     select  $S[j]$ ; arrangement
6:      $s++$ ;
7:   end if
8: end for

```



$P = (0, 1)$
I'm extracting m random numbers c_{st} forever I have in S I'm extracting a random number, according to rand number, I decide to pick or not pick

use:
• $N = 1$ want to extract only 1 item

$$\text{IF } \left(P \geq \frac{m-s}{m-j+1} \right)$$

PROB:

$$\frac{1-s}{m-j+1}$$

s starts from 0, so soon I've taken
prob is $\frac{0}{m} = 0$

S: $i_1 i_2 i_3 i_n \dots$
↑
s=0

① $\delta = 1$

$$P \leq \frac{1}{m}$$

② $P = \text{rand}(0,1)$

if $P < \frac{1-s}{m-j+1}$ \Rightarrow here prob is $\leq \frac{1}{m}$

1st item sampled with prob $\frac{1}{m}$

proof by induction (on the index j)

case not sampled i_1, i_2 to i_j

$$\begin{array}{ccccccc} i_1 & i_2 & i_3 & i_4 & \dots & i_m \\ \uparrow & \uparrow & & & & & \\ \delta = 1 & \rightarrow \delta = 2 & \delta = 3 & \vdots & & & \text{and so on} \\ P \leq \frac{1}{m} & P \leq \frac{1}{m-1} & P = \frac{1}{m-2} = \frac{1}{m-2} & \text{s still } \neq 0 \text{ (no samples)} & & & \\ & \left\{ \begin{array}{l} P \leq \frac{1-s}{m-2+1} \\ \vdots \end{array} \right. & & & & & \end{array}$$

but $\frac{1}{m}$ but is missing that if I'm on my here
 \Rightarrow means I failed sampling i_1 (before)

and so on

if I fail continuously and move to i_m

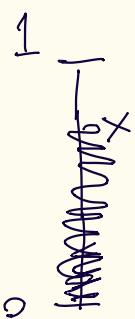
$$\delta = m$$

$$P \leq \frac{1-s}{m-m+1} = \frac{1}{1} \quad i_m \text{ is TAKEN}$$

\Rightarrow for sure
I pick 1 item

. so: I'm taking the interval $0,1$, want guarantee that algorithm takes a choice according to some given probability (our case is:
 x is a prob $0 \leq x \leq 1$)

$$P < \frac{m-s}{m-j+1}$$



$$\text{Prob}(P \leq x) = x$$

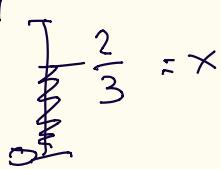
- So say if picking p uniformly ($P = \text{rand}(0,1)$)

- and compare p with x ($P < \frac{m-s}{m-j+1}$)

is equivalent to say:

this event: $P < \frac{m-s}{m-j+1}$ occurs with probability $\frac{m-s}{m-j+1} = x$

ex: if I want to execute an event with probability $\frac{2}{3}$ and I pick a random number,



$$(P \leq \frac{2}{3}) \rightarrow x$$

get p ; compare

$$p \text{ with } \frac{2}{3}$$

if It's true

\Rightarrow do the event

if not

$\Rightarrow P(\text{I execute my event}) = \text{Prob my random number } p \text{ falls here}$

but it's a uniform number

so IT's $\frac{2}{3}$ / total size = 1
 $\approx \frac{2}{3}$

Prove prob select 1 item is $\frac{1}{m}$

base case: Prob (to select the 1st item $\Rightarrow P \leq \frac{1}{m}$) = $\frac{1}{m}$

compute Prob (pick i_j)
by inductive step:
assume $P(\text{pick } i_x, \text{ for } x < j) = \frac{1}{m}$ (base case)

$$P(\text{pick } i_x, \text{ for } x < j) = \frac{1}{m}$$

$$P(i_j) = (P(\text{I don't pick any items before } i_j))$$

$$= P(\neg \text{pick } i_1, i_2, \dots, i_{j-1}) \cdot P(\text{pick } i_j)$$

$$1 - \left(\frac{P(\text{pick } i_1)}{m} + \frac{P(\text{pick } i_2)}{m} + \dots \right) =$$

$$\frac{\text{# items left to pick}}{m-j+1} = \frac{1}{m-j+1}$$

$$= \left(1 - \frac{j-1}{m} \right)$$

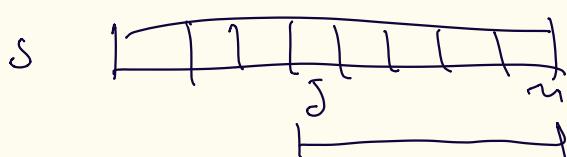
every item is picked with prob $\frac{1}{m}$

$$\Rightarrow \left(1 - \frac{j-1}{m} \right) \cdot \left(\frac{1}{m-j+1} \right) =$$

$$= \left(\frac{m-j+1}{m} \right) \cdot \left(\frac{1}{m-j+1} \right) = \frac{1}{m}$$

Generalize it to m (no pick s items but more)

$P(\text{pick } i_j \text{ given that we have to pick others } s \text{ elements})$
(generalizing in the for loop) =



missing
 $m-j+1$ items
to process

here have to pick $M-s$

items
already
picked

How many ways I can pick $m-s$ items among $m-j+1$ items?

$$\left(\text{Total # of events} = m-j+1 \text{ choose } m-s \right) = \text{Total possibility}$$

$$\binom{m-j+1}{m-s}$$

$$P = \frac{\text{Positive events} = i_j, \text{ is picked}}{\binom{m-j+1}{m-s}} = \frac{\text{count # of configurations } i_j \text{ is picked}}{\binom{m-j+1}{m-s}}$$

Picked i_j , among $m-s-1$, that can be chosen in 

$$= \frac{\binom{m-j}{m-s-1}}{\binom{m-j+1}{m-s}} = \frac{m-s}{m-j+1}$$

another proof (easy) : first notes lemma 7

 \circ

- Given the sequence $S = (a, b, c, d, e, f)$ $m=6, m=2$
 $p = (\frac{1}{2}, \frac{1}{10}, \frac{3}{10}, \frac{3}{5}, 0, 1)$

\Rightarrow extract a random number p and compare it with $\frac{m-s}{m-j+1}$

- $s=0, j=1$

$$\frac{1}{2} \stackrel{?}{<} \frac{2-0}{6-1+1} = \frac{1}{6} \quad \text{NO} \Rightarrow a \text{ is not picked}$$

$$\circ s=0, j=2 \quad \frac{1}{10} \stackrel{?}{<} \frac{2-0}{6-2+1} = \frac{2}{5} \quad \text{Yes} \Rightarrow \text{pick } b$$

$$\circ s=1, j=3 \quad \frac{3}{4} \stackrel{?}{<} \frac{2-1}{6-3+1} = \frac{1}{4} \quad \text{No}$$

$$S=1 \quad j=4 \quad \frac{3}{5} \leq \frac{2-1}{6-4+1} = \frac{1}{3} \quad \text{No}$$

$$S=1 \quad j=4 \quad 0 \leq \frac{1}{6-4+1} = \frac{1}{2} \quad \text{Yes, pick } e$$

$S=2 = b, e \Rightarrow \text{STOP}$

② N is UNKNOWN, STREAMING NODE

RESERVOIR SAMPLING

sampling with prob = $\frac{m}{n}$

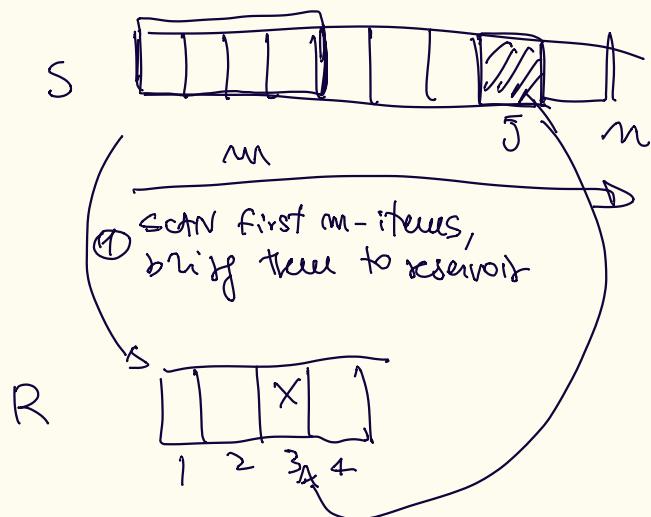
Algorithm 3.6 Reservoir sampling

```

1: Initialize array  $R[1, m] = S[1, m]$ ; # take first  $m$  items of the sequence and store in reservoir
2: for each next item  $S[j]$  do //SCAN, where  $j = m+1, \dots \Rightarrow j > m$ 
3:    $h = \text{Rand}(1, j)$ ; // 1 integer
4:   if  $h \leq m$  then
5:     set  $R[h] = S[j]$ ;
6:   end if 
7: end for
8: return array  $R$ ;

```

$m=4$



② in any position j , extract an integer(h) from 1 to j
 if the integer $(h) \leq m$
 $3 \leq 4$
 \Rightarrow substitute the item
 in pos 3 of reservoir
 and put $S[j]$

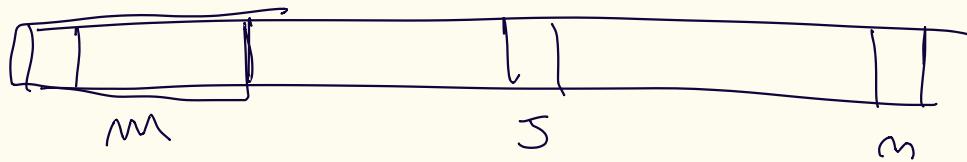
$$P(j \text{ is stored in } R) = \frac{m}{j} \quad (\text{if } h \leq m)$$

when arrive to pos. m

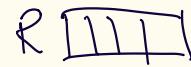
$$P(m \text{ goes in Reservoir}) = \frac{m}{m}$$

but h is taken from 1, j
 - total of possibilities $\Rightarrow j$
 - possibilities that are successful are m
 notice that $j > m$
 $(\Rightarrow P < 1)$

when I'm at the end, the probability that an item is in the reservoir is $\frac{m}{m}$



$$\frac{m}{m}$$



\Rightarrow for sure last item is in

Reservoir with $P = \frac{m}{m}$, but items before?

by induction on m : base case $m=m$ $P = \frac{m}{m} = 1$

\textcircled{R} i_j is stored in R with probability $P = \frac{m}{S}$

WANT TO SHOW

$P(\text{item } i_j \in R \text{ after } m \text{ steps}) =$
(at pos m)

i_j to be in the reservoir at the last step,
must be in R at pos $m-1$

$= P(i_j \in R \text{ after } m-1 \text{ steps}) \cdot \frac{m}{m}$

$[P(i_m \text{ is not taken}) + P(i_m \text{ is taken})]$

$$1 - \frac{m}{m}$$

$P(i_j \text{ is not picked out}) =$

$$\frac{m-1}{m} = 1 - \frac{1}{m}$$

$$= \frac{m}{m-1} \cdot \left[\left(1 - \frac{m}{m} \right) + \frac{m}{m} \cdot \frac{m-1}{m} \right] =$$

$$= \frac{m}{m-1} \cdot \left[\frac{m-m}{m} + \frac{m-1}{m} \right] =$$

$$= \frac{m}{m-1} \cdot \left[\frac{m-m+m-1}{m} \right] = \frac{m}{m-1} \cdot \frac{m-1}{m} = \frac{m}{m}$$

①



if i_j is not there, surely is not there

i_j to stay in R , two possibilities:

either i_j is not taken,

R not change
+ OR

i_m is taken

• AND i_m is not written in position of i_j

N is known, streaming model, $m=2 \quad m=8$

$$S = \{a, b, c, d, e, f, g, h\}$$

$$P = [0.5 | 0.5 | 0 | 0.5 | 1 | 0 | 1 | 1] \in \text{RAND}(0,1)$$

$s=0$

for ($j=1; j \leq n; j++$) & ($s < m$); $s++$)

$p = \text{rand}(0,1)$

if ($p < \frac{m-s}{m-j+1}$)

pick $S[j]$;

and if
else

$s++$;

$s=0 \quad j=1$

$$0.5 < \frac{m-s}{m-j+1} = \frac{2}{8-1+1} = \frac{2}{8} \cancel{>} \frac{1}{4} \Rightarrow \text{NO}$$

$s=0 \quad j=2$

$$0.5 < \frac{2}{8-2+1} = \frac{2}{7} \Rightarrow \text{NO}$$

$s=0 \quad j=3$

$$0 < \frac{2}{8-3+1} = \frac{2}{6} \cancel{>} \frac{1}{3} \Rightarrow \text{NO} \Rightarrow \text{pick } c$$

$s=1 \quad j=4$

$$\frac{1}{2} = 0.5 \stackrel{?}{<} \frac{1}{8-4+1} = \frac{1}{5} \Rightarrow \text{NO}$$

$s=1 \quad j=5$

$$0 \stackrel{?}{<} \frac{1}{8-5+1} = \frac{1}{4} \Rightarrow \text{NO}$$

$s=1 \quad j=6$

$$0 \stackrel{?}{<} \frac{1}{8-6+1} = \frac{1}{3} \text{ YES} \Rightarrow \text{pick } f$$

STOP

RESERVOIR SAMPLING (Streaming, m is unknown)

Initialize $R[1, m] = S[1, m]$

for next item $S[j] \quad // \text{item}$

$h = \text{rand}(1, j) \quad // h \text{ is an integer}$

if ($h \leq m$)

$R[h] = R[j]$

end for
end if

return R

ex: What is the status of the reservoir as soon we reach the item i?

$S = [a, b, c, d, e, f, g, h, i, \dots] \quad m=3$

2 4 1 2 3 1 : \textcircled{h} drawn from $[1, 5]$
 ↑ ↑ ↑ ↑ ↑ ↑
 current pos
of the
streaming

$R = \boxed{1 \ 2 \ 3}$

- consider that $\textcircled{1}$ is extracted in position 2 $(\frac{2}{h} \leq \frac{3}{m}) \Rightarrow R[2] = S[4]$

$R = \boxed{a \ d \ | \ c}$

- consider $\textcircled{2}$ ($\frac{4}{h} \leq \frac{3}{m}$) NO \Rightarrow Nibba

$R = \boxed{a \ d \ | \ c}$

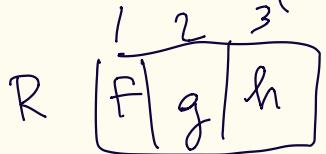
- consider $\textcircled{3}$ ($\frac{1}{h} \leq \frac{3}{m}$) yes $\Rightarrow R[1] = R[6]$

$R = \boxed{f \ | \ d \ | \ c}$

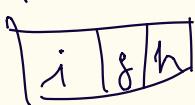
- consider $\textcircled{4}$ ($\frac{2}{h} \leq \frac{3}{m}$) yes $\Rightarrow R[2] = R[7]$

$R = \boxed{f \ | \ g \ | \ c}$

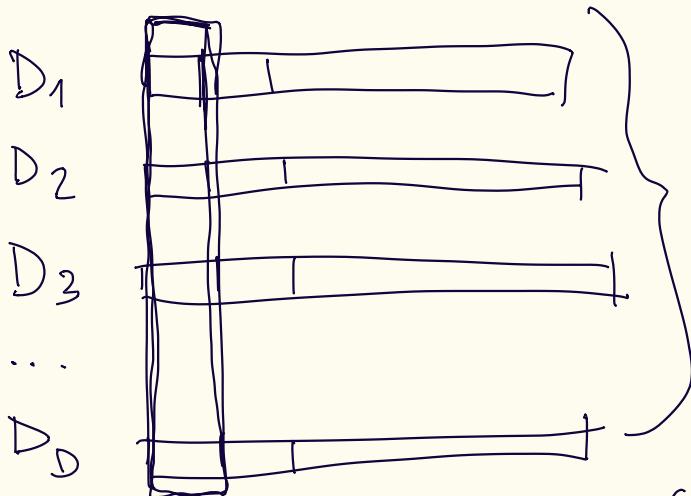
- consider (h) ($h \leq m$) yes \Rightarrow swap $R[3] = R[8]$



- consider (i) ($h \leq m$) yes \Rightarrow swap $R[1] = R[8]$



DISK STRIPING : D-disks $\Rightarrow (D > 1)$



Instead of considering the page of the disk independent, consider all pages of disks all together

1 disk, where disk page B^1 is by

$$B^1 = D \cdot B$$

In case of sorting, bound was $O\left(\frac{n}{B^1} \log_{\frac{M}{B^1}} \frac{m}{m}\right)$
over D-disk,
with disk striping

lower bound, same sort in $\log_{\frac{M}{B^1}} n$ has $\log_{\frac{M}{B^1}} \frac{n}{m}$

lower bound $\Omega\left(\frac{n}{DB}, \log_{\frac{M}{B}} \frac{m}{m}\right)$

• How much is the slow down? i.e.: how much the disk striping loses against sorting lower bound?

$$\text{RATIO} : \frac{\frac{M}{DB} \cdot \log_2 \frac{M}{DB}}{\frac{M}{B} \cdot \log_2 \frac{M}{B}} \stackrel{(\geq 1) \Rightarrow}{=} \text{wt disk striping is < then optimal}$$

$$\frac{M}{DB} \cdot \log_2 \frac{M}{DB}$$

$$\frac{\cancel{\log_2 \frac{M}{M}}}{\log_2 \frac{M}{DB}} \cdot \frac{\cancel{\log_2 \frac{M}{B}}}{\cancel{\log_2 \frac{M}{M}}} = \frac{\log_2 \frac{M}{B}}{\log_2 \frac{M}{B} - \log_2 D} = \frac{\log_2 \frac{M}{B}}{\log_2 \frac{M}{B} - \log_2 D} : \log_2 \frac{M}{B}$$

$$= \frac{1}{1 - \frac{\log_2 D}{\log_2 \frac{M}{B}}} = \frac{1}{1 - \frac{\log_2 D}{\log_2 \frac{M}{B}}} < 1 \Rightarrow \text{NFB}$$

overall > 1

The base is of thousands
and disk is tens, so

$$\Rightarrow \log_2 \frac{M}{B} D \ll 1$$

$$M \rightarrow +\infty \Rightarrow 1 \text{ OPTIMAL}$$

\Rightarrow larger is the recovery M , larger is the base,
so the log goes to 0 \Rightarrow ratio goes to 1
 \Rightarrow optimality

if $M \xrightarrow{\text{yes to}} DB \Rightarrow$ ratio goes to $\infty \Rightarrow$ slow respect to optimal algorithm

Random Sampling

$s=0$

for ($j=1 ; j \leq m \text{ } \& \text{ } s < m ; j++$)

$p = \text{rand}(0, 1)$ \Rightarrow given by the exercise

if ($p \leq \frac{m-s}{m-j+1}$)

select $S[j]$;

~~$s += 1$~~ ;

and if and for

$$p \leq \frac{m-s}{m-j+1}$$

Random Sampling

$m=2$ $m=8$ items $S = \{a, b, c, d, e, f, g, h\}$

$j=1 \quad s=0$

$$p = \frac{1}{2}$$

$$p \leq \frac{2-0}{8-1+1} = \frac{2}{8} = \frac{1}{4} \text{ NOT picked}$$

$j=2 \quad s=0$

$$p = \frac{1}{2} \leq \frac{2-0}{8-2+1} = \frac{2}{7} \text{ NOT picked}$$

$j=3, s=0$

$$p = 0 \leq \frac{2-0}{8-3+1} = \frac{2}{6} = \frac{1}{3} \text{ yes, } \boxed{C} \text{ is picked}$$

$j=4; s=1$

$$p = \frac{1}{2} \leq \frac{2-1}{8-4+1} = \frac{1}{5} \text{ NOT picked}$$

$j=5; s=1$

$$p = 1 \leq \frac{2-1}{8-5+1} = \frac{1}{4} \text{ NOT picked}$$

$j=6; s=1$

$$p = 0 \leq \frac{2-1}{8-6+1} = \frac{1}{3} \text{ yes, } \boxed{F} \text{ is picked}$$

$s=2 \quad I \quad \text{STOP}$

RESERVOIR SAMPLING

$$m=2 \quad S = \{ \underset{1}{a}, \underset{2}{b}, \underset{1}{c}, \underset{3}{d}, \underset{2}{e}, \underset{1}{f} \}$$

$$R = \begin{array}{|c|c|} \hline 1 & 2 \\ \hline a & b \\ \hline \end{array}$$

↑ ↑ ↑ ↑

$h(c) = 1 \leq m=2$ yes $\Rightarrow R[1] = c$

$$R = \begin{array}{|c|c|} \hline 1 & 2 \\ \hline c & b \\ \hline \end{array}$$

$h(d) = 3 \geq m=2$ NO, NisBA

(still $R: \boxed{c | b}$)

$h(e) = 2 \leq m=2$ yes $\Rightarrow R[2] = e$

$$R = \begin{array}{|c|c|} \hline 1 & 2 \\ \hline c & e \\ \hline \end{array}$$

$h(f) = 1 \leq m=2$ yes $\Rightarrow R[1] = f$

$$R = \begin{array}{|c|c|} \hline 1 & 2 \\ \hline f & e \\ \hline \end{array}$$
✓

MULTI-KEY QS

pivot always first step
in n

R: (BUS, BATH, ABACUS, AARGH, CAT)

i=1

?

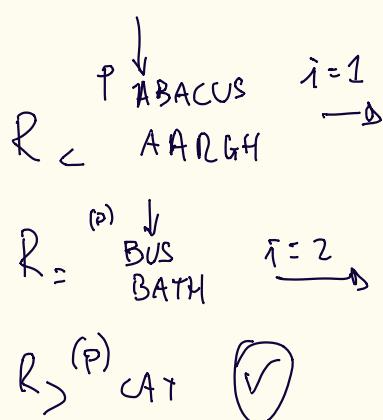
(BUS)

BATH

ABACUS

AARGH

CAT



R < $\min(R, P)$

R = P

AARGH

i=2

R > -

R < $\min(R, P)$

R = P

BATH

i=2

R > -

R < $\min(R, P)$

R = P

BUS

i=3

R > -

R < AARGH

i=2

R = ABACUS

i=3

R >



Reservoir Sampling

$m = 2$

$$S = \begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 \\ a & b & c & d & e & f \\ 3 & 1 & 4 & 2 \end{pmatrix}$$

$$R = \boxed{\begin{array}{c|c} 1 & 2 \\ \hline A & B \end{array}}$$

$$h(c) \stackrel{?}{\leq} m \quad \text{no}$$

$$h(d) \stackrel{?}{\leq} 2 \quad \text{yes}$$

$$R[2] = S[4]$$

$$R = \boxed{\begin{array}{c|c} 1 & 2 \\ \hline D & B \end{array}}$$

$$h(e) \stackrel{?}{\leq} 2 \quad \text{no}$$

$$h(f) \stackrel{?}{\leq} 2 \quad \text{yes} \quad R[2] = S[6]$$

$$R = \boxed{\begin{array}{c|c} 1 & 2 \\ \hline D & F \end{array}}$$