

## Example 2

② sketch: naive algorithm

Given a set of  $N$  strings, sample a random subset of size  $M$ , print sorted. (arbitrary to input)

all set is available in memory (classic)  
every string seen just once (streaming)

obs: Given two integers  $N, M$ , where  $M \leq N$ , sample a random subset of size  $M$  from  $[1, N]$ , output sorted.

Tools:  $\text{RandInt}(I, J)$ : generates random integer in  $[I, J]$   
 $\text{RandReal}(0, 1)$ : real in  $[0, 1)$

① Generate  $M$  random integers  $\text{RandInt}(1, N)$ , then sort

- Time:  $O(M \lg M)$  mergesort

- space:  $O(M)$

- Cons: duplicates? possibly generate less. 

no time optimal, because of "log-slowdown".

① set a bit in an array  $\Rightarrow$  no sort, yes scan, but  $O(N)$  space

② Generate w/outly  $M$  ints by SCAN  $\Rightarrow$  no sort and  $M$  space

$\rightarrow$  useful in a streaming environment

$\rightarrow$  Pick item  $i$  with  $P = \frac{\# \text{remaining elements}}{\# \text{items left}}$

$m$  selected

  $P(\text{pick } i) = \frac{M - m}{N - i}$

- Time =  $O(N)$  [may be more than  $M \lg M$ , but scan-based]

- space =  $M$  [already sorted, output in order when select]

CONS if  $M \ll N$

this is bad

[think to select  $m$  edges out of  $M$  in a graph  
 $\parallel$  MST extension]

③

~~Generate "gaps" rather than "indicators"~~

Deploy randomness, good "average case"

- ~~M~~ M bins of size  $N/M$ : integer  $i$  goes to  $i \times \frac{M}{N}$
- bins implemented as lists.
- because ints are random, each list has average length 1 (hint: bucket sort)

- Time:  $O(M)$  average (even insertion sort is ok on each bin)

- Space:  $M$

CONS: resampling, average complexity  $\rightarrow$   $\left\{ \begin{array}{l} m \text{ extractions have prob of} \\ \text{a conflict} \leq \frac{1}{N} + \frac{2}{N} + \dots + \frac{m-1}{N} = \frac{m^2}{N} \end{array} \right.$

④

Generate "gaps" rather than "indicators".

$G(\text{gap}, n) \Rightarrow$  random var that indicates the gap of the next element, where  $m$  have to be still selected out of  $n$  left.

$$P(G=g) = f(g) = \frac{\binom{n-g-1}{m-1}}{\binom{n}{m}}$$

- The key problem is how to generate  $g$  according to  $f()$ .
- Vitter proposed a sophisticated solution for this
  - paper attached
  - linear time (on average)
  - no conflicts and ~~no~~ no resampling

SEQUENCE TO SHOW

	time	space	features
1'	$N$	$N$	scan only (types), no sort
2	$N$	$M$	knuth (streaming) NO-RESAMPLE
1	$M \log M$	$M$	softmax + resample
3	$M$	$M$	average + resample <span style="border: 1px solid red; padding: 2px;">VITTER ④</span> + stream