
L7: Reservoir Sampling

scribe(s): *Kavitaben Sheth*

Streaming algorithms are process that takes data in input stream and examine in few passes typically one pass. These algorithms has limited memory available and limited processing time to process each data item. All sampling methods that process the input file in one pass can be characterised as reservoir algorithms.

Reservoir Sampling is a family of randomized but fast algorithms for selecting a random sample of n records without replacement from a pool of N records, where value of N is unknown beforehand. Reservoir algorithm select first n records of the file into a 'reservoir' and rest of records are process sequentially.

7.1 Reservoir algorithm

Random sampling is basic to many computer applications in computer science, statistics, and engineering. The problem is to select without replacement a random sample of size n from a set of size N .

The first step of any reservoir algorithm is to put the first n records of the file into a reservoir. The rest of the records are processed sequentially; records can be selected for the reservoir only as they are processed. An algorithm is a reservoir algorithm if it maintains the invariant that after each record is processed a true random sample of size n can be extracted from the current state of the reservoir. This algorithm runs in $O(N)$ because the entire file must be scanned.

So, we assume we have a continuous random variable that can be generated quickly and whose distribution approximate well. It can be generated very quickly, in constant time. So, At the i th element of S , the algorithm generates a random number j between 1 and i . If j is less than n , the j th element of the reservoir array is replaced with the i th element of S . In effect, for all i , the i th element of S is chosen to be included in the reservoir with probability n/i . Assume $i+1$ to m elements kept in reservoir. Now, probability of not replaced element from reservoir is $1-(m-i)/m = i/m$. Total probability of kept and not replaced elements is $(n/i)*(i/m) = n/m$. So, when algorithm finished executing, each element has equal probability of being chosen from reservoir.

7.2 Approximate Counts

The approximate counting problem counts large number of events using small amount of memory.

Consider query Interval I subset with elements $[n]$ $count(I) = |a_i \text{ in } A |_{a_i \text{ in } I}$

Goal is to find data structure S such that for every query interval probability = $|S(I) - count(I)| > \epsilon * m$. Probability $< \delta$.

7.2.1 Chernoff Inequality

Let X_1, X_2, \dots, X_r be independent reservoirs.

Let $\Delta_i = \max(X_i) - \min(X_i)$

Let $M = \sum_{i=1}^r X_i$

$$\text{Probability} = Pr[|M - \sum_{i=1}^r E[X_i]| \geq \alpha] \leq 2 \exp\left(\frac{-2\alpha^2}{\sum_i \Delta_i^2}\right)$$

Often : $\Delta = \max_i \Delta_i$ and $E[X_i] = 0$ then :

$$\text{Probability} = Pr[|M| \geq \alpha] \leq 2 \exp\left(\frac{-2\alpha^2}{r\Delta^2}\right)$$

Let S be a random sample of size $k = O\left(\frac{1}{\epsilon^2} \log\left(\frac{1}{\Delta}\right)\right)$

$$S(I) = |Scap I| * \left(\frac{m}{k}\right)$$

Each s_i in I with probability $(\text{count}(I)/m)$

Reservoir $RV Y_i = \{1 \text{ if } s_i \text{ in } I, 0 \text{ if } s_i \text{ not in } I\}$

$$E[Y_i] = \text{count}(I)/m$$

Reservoir $RV X_i = (Y_i - \text{count}(I)/m)/k$

$$E[X_i] = 0$$

$$\Delta < 1/k$$

$M = \text{sum}_i X_i$ == error on count estimate by S

$$\Pr[|M| > \epsilon] < 2 \exp\left(\frac{-2\epsilon^2}{k * (1/k^2)}\right) < \Delta$$