

Reservoir Sampling

Reservoir sampling is a family of randomized algorithms for randomly choosing k samples from a list of n items, where n is either a very large or unknown number. Typically n is large enough that the list doesn't fit into main memory. For example, a list of search queries in Google and Facebook.

So we are given a big array (or stream) of numbers (to simplify), and we need to write an efficient function to randomly select k numbers where $1 \leq k \leq n$. Let the input array be *stream[]*.

A **simple solution** is to create an array *reservoir[]* of maximum size k . One by one randomly select an item from *stream[0..n-1]*. If the selected item is not previously selected, then put it in *reservoir[]*. To check if an item is previously selected or not, we need to search the item in *reservoir[]*. The time complexity of this algorithm will be $O(k^2)$. This can be costly if k is big. Also, this is not efficient if the input is in the form of a stream.

It **can be solved in $O(n)$ time**. The solution also suits well for input in the form of stream. The idea is similar to [this](#) post. Following are the steps.

- 1) Create an array *reservoir[0..k-1]* and copy first k items of *stream[]* to it.
- 2) Now one by one consider all items from $(k+1)$ th item to n th item.
 - ...a) Generate a random number from 0 to i where i is index of current item in *stream[]*. Let the generated random number is j .
 - ...b) If j is in range 0 to $k-1$, replace *reservoir[j]* with *arr[i]*

Following is C implementation of the above algorithm.

```
// An efficient program to randomly select k items from :

#include <stdio.h>
#include <stdlib.h>
#include <time.h>

// A utility function to print an array
void printArray(int stream[], int n)
{
```

```
// A function to randomly select k items from stream[0..n-1]
void selectKItems(int stream[], int n, int k)
{
    int i; // index for elements in stream[]

    // reservoir[] is the output array. Initialize it with
    // first k elements from stream[]
    int reservoir[k];
    for (i = 0; i < k; i++)
        reservoir[i] = stream[i];

    // Use a different seed value so that we don't get
    // same result each time we run this program
    srand(time(NULL));

    // Iterate from the (k+1)th element to nth element
    for (; i < n; i++)
    {
        // Pick a random index from 0 to i.
        int j = rand() % (i+1);

        // If the randomly picked index is smaller than
        // the element present at the index with new element
        if (j < k)
            reservoir[j] = stream[i];
    }

    printf("Following are k randomly selected items \n");
    printArray(reservoir, k);
}
```

data:text/html;charset=utf-8,%3Cdiv%20class%3D%22post-title-info%22%20style%3D%22float%3A%20left%3B%20font-size%... 2/4

Output:

Following are k randomly selected items
6 2 11 8 12

Time Complexity: $O(n)$

How does this work?

To prove that this solution works perfectly, we must prove that the probability that any item $stream[i]$ where $0 \leq i < n$ will be in final $reservoir[]$ is k/n . Let us divide the proof in two cases as first k items are treated differently.

Case 1: For last $n-k$ stream items, i.e., for $stream[i]$ where $k \leq i < n$

For every such stream item $stream[i]$, we pick a random index from 0 to i and if the picked index is one of the first k indexes, we replace the element at picked index with $stream[i]$

To simplify the proof, let us first consider the *last item*. The probability that the last item is in final reservoir = The probability that one of the first k indexes is picked for last item = k/n (the probability of picking one of the k items from a list of size n)

Let us now consider the *second last item*. The probability that the second last item is in final $reservoir[]$ = [Probability that one of the first k indexes is picked in iteration for $stream[n-2]$] X [Probability that the index picked in iteration for $stream[n-1]$ is not same as index picked for $stream[n-2]$] = $[k/(n-1)] * [(n-1)/n] = k/n$.

Similarly, we can consider other items for all stream items from $stream[n-1]$ to $stream[k]$ and generalize the proof.

Case 2: For first k stream items, i.e., for $stream[i]$ where $0 \leq i < k$

The first k items are initially copied to $reservoir[]$ and may be removed later in iterations for $stream[k]$ to $stream[n]$.

The probability that an item from $stream[0..k-1]$ is in final array = Probability that the item is not picked when items $stream[k]$, $stream[k+1]$, $stream[n-$

1] are considered = $\frac{k}{k+1} \times \frac{k+1}{k+2} \times \frac{k+2}{k+3} \times \dots \times \frac{n-1}{n} = \frac{k}{n}$

References:

http://en.wikipedia.org/wiki/Reservoir_sampling