

Random Sampling

Florian Schoppmann

August 24, 2010

Sampling Algorithms

Input:

- List[1 . . . N]
- Length of database N (if known)
- Length of sample n

Output:

- Sample[1 . . . n]

Considerations

- Online vs. random-access
- Sequential vs. non-sequential
- Samples for independent categories

Desiderata:

- Parallelizable
- If random access, running time close to $O(n)$
- Constant memory

Random Indices

```
1:  $m \leftarrow 0$ 
2: while  $m < n$  do
3:    $R \leftarrow \text{random}(\{1 \dots N\})$ 
4:   if  $\text{List}[R] \notin \text{Sample}$  then
5:      $m \leftarrow m + 1$ 
6:      $\text{Sample}[m] \leftarrow \text{List}[R]$ 
```

- ca. $N \ln \frac{N}{N-n+1}$ iterations in expectation
- Space/time trade off in line 4

Random Remaining Indices

- 1: **for** $m \leftarrow 1, \dots, n$ **do**
- 2: $R \leftarrow \text{random}(\{1 \dots N - m + 1\})$
- 3: $j \leftarrow$ index of R 'th non-null element in List
- 4: $\text{Sample}[m] \leftarrow \text{List}[j]$
- 5: $\text{List}[j] \leftarrow \text{null}$

- Prohibitive running time $\Theta(nN)$
- Modifies List

The Fisher-Yates Shuffle

```
1: for  $m \leftarrow 1, \dots, n$  do  
2:    $R \leftarrow \text{random}(\{m \dots N\})$   
3:   Swap List[ $m$ ] and List[ $R$ ]  
4: Sample[ $1 \dots n$ ]  $\leftarrow$  List[ $1 \dots n$ ]
```

- Running time $\Theta(n)$
- Modifies List

Probabilistic Sampling

```
1: for  $t \leftarrow 1, \dots, N$  do  
2:   with probability  $\frac{n}{N}$  do  
3:     Append List[ $t$ ] to Sample
```

- Running time $\Theta(N)$
- Only expected sample size n (mean of $B(N, \frac{n}{N})$)
- Standard deviation $\sqrt{n(1 - n/N)}$

Selection Sampling

```
1:  $m \leftarrow 0$ 
2: for  $t \leftarrow 1, \dots, N$  do
3:     with probability  $\frac{n-m}{N-t}$  do
4:          $m \leftarrow m + 1$ 
5:          $\text{Sample}[m] \leftarrow \text{List}[t]$ 
```

- Running time $\Theta(N)$
- Completely unbiased!

Random Number Generation

(Digression)

Running time $O(n)$ possible by skipping rows?

Idea 1:

- Let $S \in \{0 \dots N - n\}$ RV for # rows to skip

$$\Pr[S \leq s] = 1 - \frac{(N - n)^{s+1}}{N^{s+1}}$$

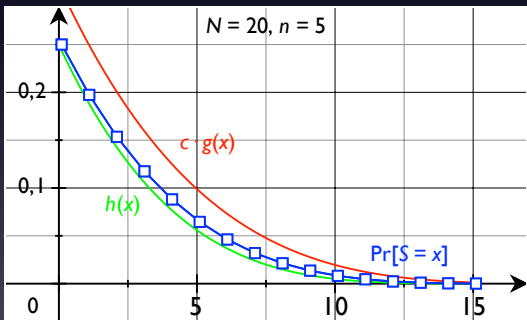
Idea 2 (Vitter, 1984):

- von Neumann's rejection & "squeeze" method

Vitter (1984)

$$c = \frac{N}{N-n+1}, \quad g(x) = \frac{n}{N} \left(1 - \frac{x}{N}\right)^{n-1},$$

$$h(x) = \frac{n}{N} \left(1 - \frac{x}{N-n+1}\right)^{n-1}$$



Reservoir Sampling

```
1: Sample[1 . . . n] ← List[1 . . . n]
2: for  $t \leftarrow n + 1, \dots, N$  do
3:   with probability  $\frac{n}{t}$  do
4:      $R \leftarrow \text{random}(\{1 \dots n\})$ 
5:     Sample[R] ← List[t]
```

- Completely unbiased!
- $O(n(1 + \log \frac{N}{n}))$ by optimizing (Vitter, 1985)

Reservoir, with Replacement

```
1: for  $t \leftarrow 1, \dots, N$  do  
2:   for  $i \leftarrow 1, \dots, n$  do  
3:     with probability  $\frac{1}{t}$  do  
4:        $\text{Sample}[i] \leftarrow \text{List}[t]$ 
```

- Completely unbiased!

Bibliography



Knuth (1997):

The Art of Computer Programming, Vol. 2



Vitter (1984):

Faster methods for random sampling



Vitter (1985):

Random sampling with a reservoir



Park, Ostrouchov, Samatova, Geist (2004):

Reservoir-Based Random Sampling with
Replacement from Data Stream