

Sampling a Stream

What Doesn't Work

Sampling Based on Hash Values

When Sampling Doesn't Work

- Suppose Google would like to examine its stream of search queries for the past month to find out what fraction of them were unique – asked only once.
- But to save time, we are only going to sample $1/10^{\text{th}}$ of the stream.
- The fraction of unique queries in the sample will not be the fraction for the stream as a whole.
 - In fact, we can't even adjust the sample's fraction to give the correct answer.

Example: Unique Search Queries

- The length of the sample is 10% of the length of the whole stream.
- Suppose a query is unique.
 - It has a 10% chance of being in the sample.
- Suppose a query occurs exactly twice in the stream.
 - It has an 18% chance of appearing exactly once in the sample.
- And so on ... The fraction of unique queries in the stream is unpredictably large.

Sampling by Value

- **Our mistake:** we sampled based on the position in the stream, rather than the value of the stream element.
- Hash search queries to 10 buckets 0, 1,..., 9.
- Sample = all search queries that hash to bucket 0.
 - All or none of the instances of a query are selected.
 - Therefore the fraction of unique queries in the sample is the same as for the stream as a whole.

Controlling the Sample Size

- **Problem:** What if the total sample size is limited?
- **Solution:** Hash to a large number of buckets.
- Adjust the set of buckets accepted for the sample, so your sample size stays within bounds.

Example: Fixed Sample Size

- Suppose we start our search-query sample at 10%, but we want to limit the size.
- Hash to, say, 100 buckets, 0, 1,..., 99.
 - Take for the sample those elements hashing to buckets 0 through 9.
- If the sample gets too big, get rid of bucket 9.
- Still too big, get rid of 8, and so on.

Sampling Key-Value Pairs

- Our technique generalizes to any form of data that we can see as tuples (K, V) , where K is the “key” and V is a “value.”
- **Distinction**: We want our sample to be based on picking some set of keys only, not pairs.
 - In the search-query example, the data was “all key.”
- Hash keys to some number of buckets.
- Sample then consists of all key-value pairs with a key that goes into one of the selected buckets.

Example: Salary Ranges

- Data = tuples of the form (EmpID, Dept, Salary).
- **Query**: What is the average range of salaries within a department?
- Key = Dept.
- Value = (EmpID, Salary).
- Sample picks some departments, has salaries for all employees of that department, including its min and max salaries.