**Note to other teachers and users of these slides:** We would be delighted if you found this our material useful in giving your own lectures. Feel free to use these slides verbatim, or to modify them to fit your own needs. If you make use of a significant portion of these slides in your own lecture, please include this message, or a link to our web site: <u>http://www.mmds.org</u>

# Sampling from a stream

Adapting material from: Mining of Massive Datasets Jure Leskovec, Anand Rajaraman, Jeff Ullman Stanford University http://www.mmds.org



## Data Streams

In many data mining situations, we do not know the entire data set in advance

- Stream Management is important when the input rate is controlled externally:
  - Google queries
  - Twitter or Facebook status updates
- We can think of the data as infinite and non-stationary (the distribution changes over time)

## **The Stream Model**

- Input elements enter at a rapid rate, at one or more input ports (i.e., streams)
  We call elements of the stream tuples
- The system cannot store the entire stream accessibly
- Q: How do you make critical calculations about the stream using a limited amount of (secondary) memory?

#### Sampling from a Data Stream: Sampling a fixed proportion

As the stream grows the sample also gets bigger

#### Sampling from a Data Stream

- Since we can not store the entire stream, one obvious approach is to store a sample
  Two different problems:
  - (1) Sample a fixed proportion of elements in the stream (say 1 in 10)
  - (2) Maintain a random sample of fixed size over a potentially infinite stream
    - At any "time" k we would like a random sample of s elements
      - What is the property of the sample we want to maintain? For all time steps k, each of k elements seen so far has equal prob. of being sampled

## Sampling a Fixed Proportion

- Problem 1: Sampling fixed proportion
- Scenario: Search engine query stream
  - Stream of tuples: (user, query, time)
  - Answer questions such as: How often did a user run the same query in a single days
  - Have space to store 1/10<sup>th</sup> of query stream

#### Naïve solution:

- Generate a random integer in [0..9] for each query
- Store the query if the integer is 0, otherwise discard

## Problem with Naïve Approach

- Simple question: What fraction of queries by an average search engine user are duplicates?
  - Suppose each user issues x queries once and d queries twice (total of x+2d queries)
    - Correct answer: d/(x+d)

#### Proposed solution: We keep 10% of the queries

- Sample will contain x/10 of the singleton queries and 2d/10 of the duplicate queries at least once
- But only *d*/100 pairs of duplicates
  - d/100 = 1/10 · 1/10 · d
- Of d "duplicates" 18d/100 appear exactly once
  - 18d/100 = ((1/10 · 9/10)+(9/10 · 1/10)) · d



# **Solution: Sample Users**

#### **Solution:**

- Pick 1/10<sup>th</sup> of users and take all their searches in the sample
- Use a hash function that hashes the user name or user id uniformly into 10 buckets

# **Generalized Solution**

#### Stream of tuples with keys:

- Key is some subset of each tuple's components
  - e.g., tuple is (user, search, time); key is user
- Choice of key depends on application

#### To get a sample of *a*/*b* fraction of the stream:

- Hash each tuple's key uniformly into b buckets
- Pick the tuple if its hash value is at most a



Hash table with **b** buckets, pick the tuple if its hash value is at most **a**. **How to generate a 30% sample?** Hash into b=10 buckets, take the tuple if it hashes to one of the first 3 buckets