

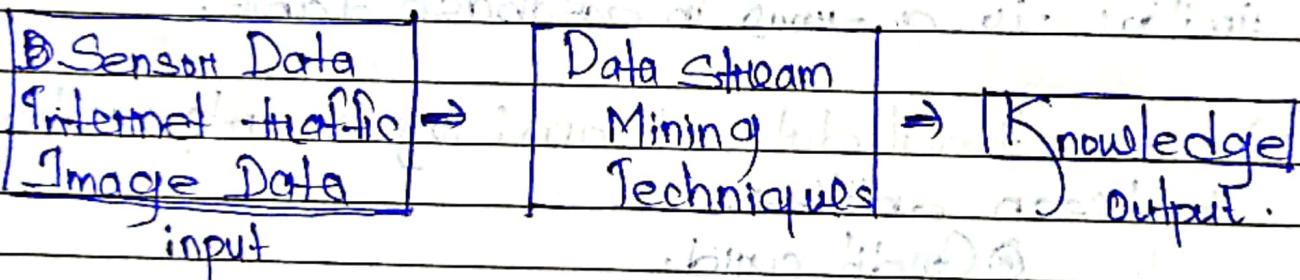
- Data Stream: Continuous, fast-changing, and ordered chain of data - transmitted at a very high speed.
- ⇒ An ordered sequence of information for a specific interval.
- ⇒ Sender's data is transferred from sender's side and immediately shows in data streaming at the receiver's side.
- ⇒ Streaming does not mean downloading the data or storing the information on storage devices.
- ⇒ A data stream is potentially unbounded sequence of tuples. Each tuple consists of a set of attributes similar to a row in database table.
- ⇒ Transactional data stream: log of interactions between entities.
 - ④ Credit card.
 - ④ Telecommunications
 - ④ Web
- ⇒ Measurement data streams: monitor evolution of entity states
 - ④ Sensor networks = road traffic
 - ④ IP networks = traffic at router interfaces
 - ④ Earth climate = temperature, moisture

Sources of data stream:

- i) Internal traffic
- ii) Sensors data
- iii) Real-time ATM transaction
- iv) Live event data.
- v) Call records
- vi) Satellite data
- vii) Audio listening
- viii) Watching videos.

Data stream

Generators



Sensor Data: In navigation systems, sensor data is used. Imagine a temperature sensor floating about in the ocean, sending back to the base station a reading of the surface temperature each hour. The data generated by this sensor is a stream of real numbers.

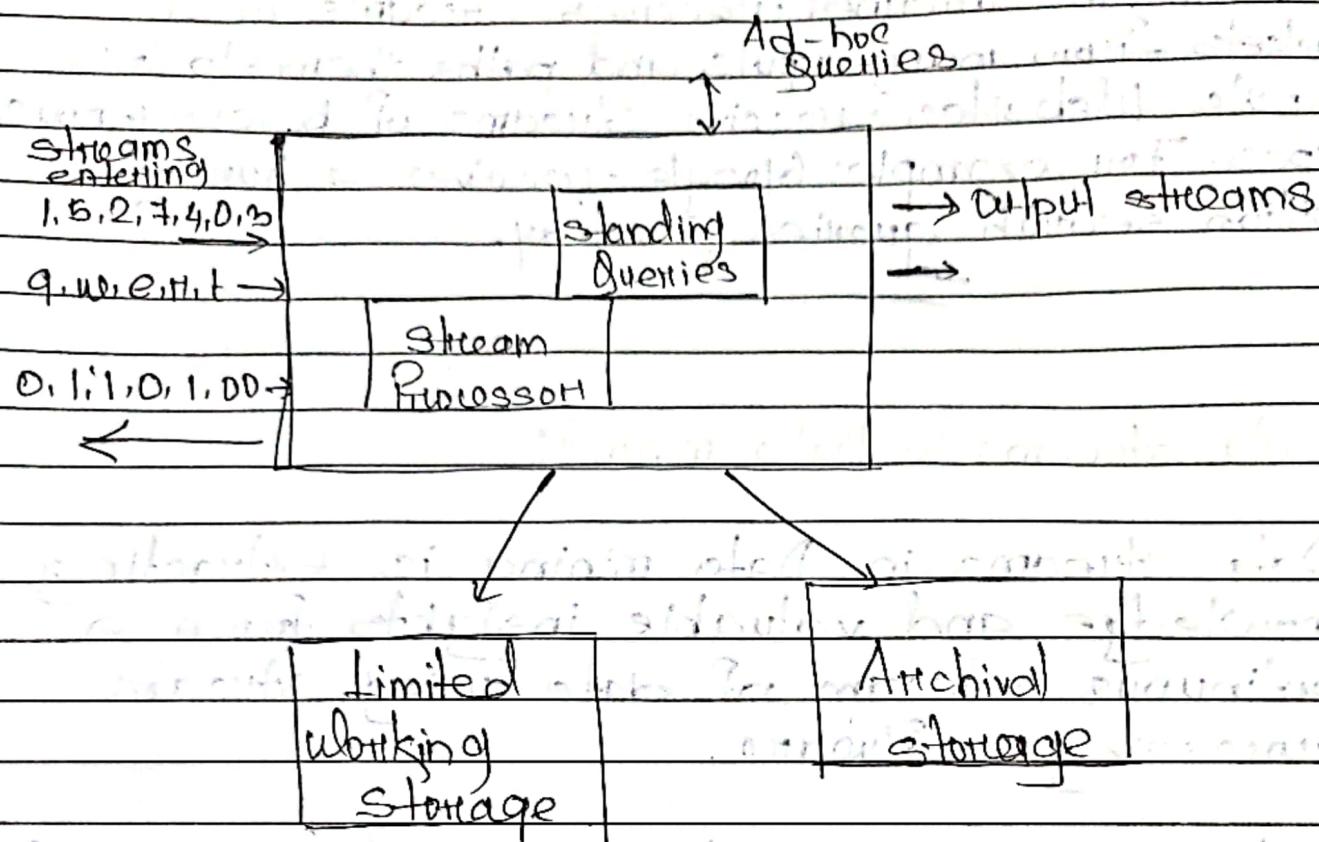
Image Data: Satellites frequently send down to each earth streams containing many terabytes of images per day. Surveillance cameras generate images with lower resolution than satellites, but there can be numerous of them, each producing a stream of images at a rate of 1 second.

ii) Internet and Web Traffic: A bobbing node in the center of the internet receives streams of IP packets from many inputs and paths them to its outputs. Websites receive streams of heterogeneous types. For example: Google receives a hundred million search queries per day.

Data streams in Data mining:

- ⇒ Data streams in Data mining is extracting knowledge and valuable insights from a continuous stream of data using stream processing software.
- ⇒ Data streams in data mining can be considered a subset of general concepts of machine learning, knowledge extraction and data mining.
- ⇒ Data analysis of large amount of data needs to be done in real-time.
- ⇒ The structure of knowledge is extracted is represented in the case of models and patterns of infinite streams of information.

A data-stream management system (DSMS):



→ Streams may be archived in a large archival store, but we assume it is not possible to answer queries from the archival store.

→ DSMS could be examined only under special circumstances using time-consuming retrieval process.

→ There is also a working store, into which summaries on parts of streams may be placed, and which can be used for answering queries.

→ The working store might be disk, or it might be main memory, depending on how fast we need to process queries.

- But either way, it is of sufficiently limited capacity that it cannot store all the data from all the streams.
- Any number of streams can enter the system.
- Each stream can provide elements at its own schedule; they need not have the same data rates or data types, and the times between elements of one stream need not be uniform.

■ Problems on Data Streams:

- i) Sampling data from a stream: Construct a random sample.
- ii) Queries over sliding windows: Number of items of type x in the last k elements of the stream.
- iii) Filtering a data stream: Select elements with property x from the stream.
- iv) Counting distinct elements: Number of distinct elements in the last k elements of the stream.
- v) Estimating moments: Estimate avg / std dev of last k elements.
- vi) Finding frequent elements.

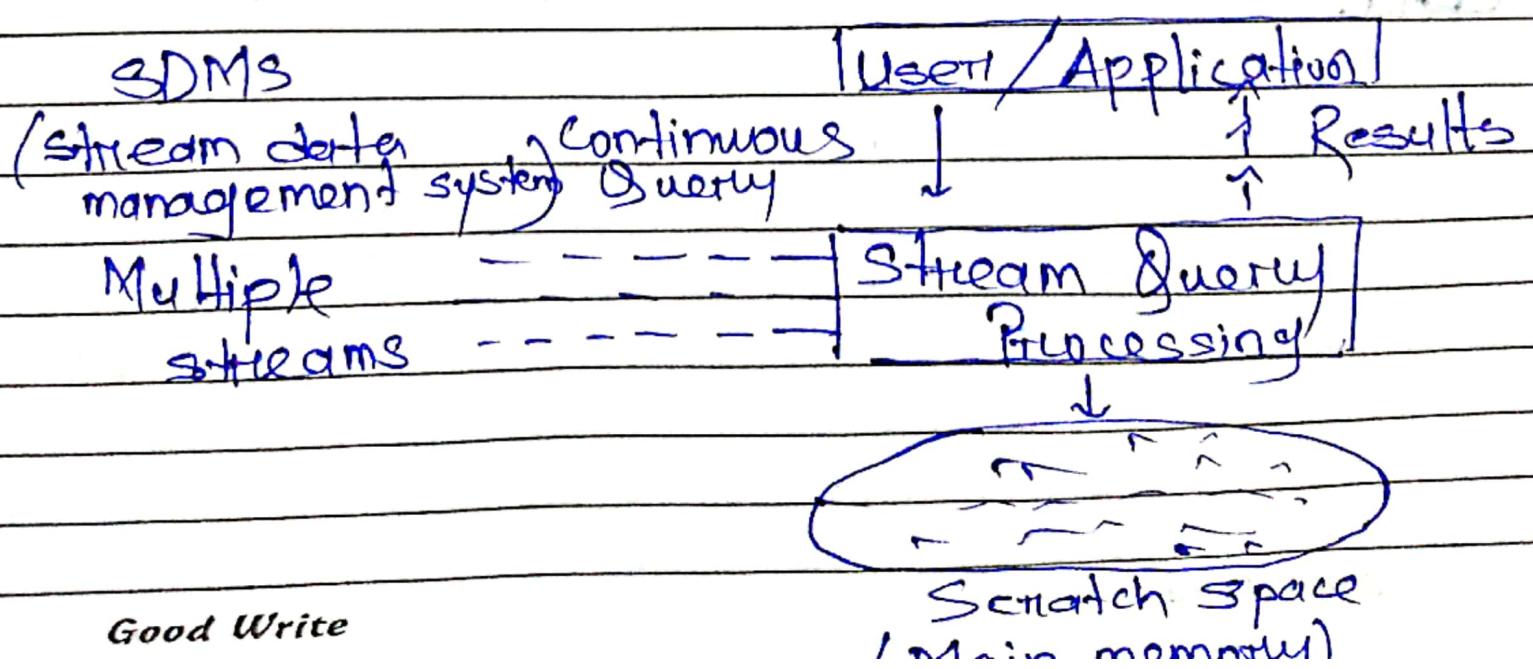
Applications:

- i) Mining query streams: Google wants to know what queries are more frequent today than yesterday.
- ii) Mining click streams: Yahoo wants to know which of its pages are getting an unusual number of hits in the past hour.
- iii) Mining social network newsfeed: look for trending topics on Twitter, Facebook.
- iv) Sensor networks: Many sensors feeding into a central controller.
- v) Telephone call records: Gather information for optimal routing, Detect denial-of-service attacks.
- vi) Data stream processing: Process queries, apply data mining algorithms.
- vii) Requirements:
 - Real-time processing
 - One-pass processing
 - Bounded storage.
 - Possibly consider several streams.

Characteristics of Data Stream in data mining:

- i) Continuous stream of data: Infinite continuous stream resulting in data big data. In data streaming, multiple data streams are passed simultaneously.
- ii) Time sensitive: Elements of data streams carry time stamps with them. After a particular time, the data stream loses its significance.
- iii) Data Volatility: No data is stored in data streaming. Once the data mining and analysis are done, information is summarized or discarded.
- iv) Concept Drifting: Very unpredictable. The data changes or evolves with time; as in this dynamic world, nothing is constant.

Stream Query Processing architecture:



Difference between Database Management System and data stream management system:

DBMS

- i) Permanent updatable relations.
- ii) Data stored in disks.
- iii) SQL language.
- iv) Performs in large volume of data.
- v) One-time queries.
- vi) Random access.
- vii) Unbounded disk storage.
- viii) Only current state matters.
- ix) No real time service.
- x) Assume precise data.

DSMS

- i) Streams and permanent updatable relations.
- ii) Streams are processed on fly and permanent relation stored in disk.
- iii) SQL like query language.
- iv) Optimization of computer resources to deal with several streams, several queries.
- v) Continuous queries.
- vi) Sequential access.
- vii) Bounded disk storage.
- viii) Historical data is important.
- ix) Real time requirements.
- x) Data state/imprecise.

Challenges of stream data processing:

- i) Multiple, continuous, rapid, time-varying, ordered streams.
- ii) Main memory computations.
- iii) Queries are often continuous.
- iv) Queries are often complex.
- v) Multi-level / Multi-dimensional processing of data mining.

Stamping Computing Approaches:

⇒ Two approaches for handling such streams:

- i) Use a time window
- ii) Query the window as static table

⇒ When one can't store collected data, or

to keep tracks of historical data

- i) Sampling
- ii) Filtering
- iii) Counting

Data streams in data mining techniques:

⇒ Data streams in data mining techniques are implemented to extract patterns and insights from a data stream.

Data Stream Mining Techniques:

- i) Classification: Classification is a supervised learning technique. In classification, the classifier model is built based on the training data.
- This classifier model is then used to predict the label for unlabeled instances or items continuously arriving through the data stream.
 - Prediction is made for the unknown/new items that the model never seen, and already known instances are used to train the model.
 - Receives an unlabeled item and predict it based on its current model.
 - Receives labels for past known items and use them for the training the model.
- Classification techniques:
- i) Lazy Classifier: or k-nearest Neighbour.
 - ii) Naïve Bayes
 - iii) Decision Tree
 - iv) Logistic Regression.

- ii) Regression: Supervised learning technique used to predict real values of label attributes for the stream instances, not the discrete values like classification.
- The idea of regression is similar to classification either to predict the real values label for unknown items using the regression model.

Algorithms:

- i) Lazy classifier on K-nearest Neighbour
 - ii) Naive Bayes
 - iii) Decision Trees
 - iv) Linear Regression
 - v) Ensembles.
- iii) Clustering: Unsupervised learning technique.

- ⇒ Clustering is functional when we have unlabeled instances, and we want to find homogenous clusters in them based on the similarities of data item.
- ⇒ Before the clustering process, the groups are not known. Clusters are formed with continuous data streams based on date and keep on adding items to the different groups.

- Algorithms:
- i) K-means clustering
 - ii) Hierarchical Clustering
 - iii) Density-based clustering

w) Frequent Pattern mining: An essential task in unsupervised learning.

⇒ It is used to describe the data and find the association rules or discriminative features in data that will further help classification and clustering tasks. Based on two rules:

- i) Frequent Item Set: Collection of items occurring together frequently.
- ii) Association Rules: Indicator of the strong relationship between two items

Sampling from a Data Stream:

Since we cannot store the entire stream, one obvious approach is to store a sample.

Two different problems:

- i) Sample a fixed proportion of elements in the stream.
- ii) Maintain a random sample of fixed size over a potentially infinite stream.

i) Problem 1: Sampling fixed proportion.

Scenario \Rightarrow Search engine query stream
 \rightarrow Stream of tuples: (user, query, time)
 \rightarrow Answer questions such as:
- "How often did a user run the same query in a single day?"
- Have space to store $1/10^{\text{th}}$ of query stream.

Naive solution \Rightarrow

- i) Generate a random integer in $[0 \dots 9]$ for each query.
- ii) Store the query if the integer is 0, otherwise discard.

Solution \Rightarrow

- i) Pick $1/10^{\text{th}}$ user's and take all their searches in the sample.
- ii) Use a hash function that hashes the user name or user id uniformly into 10 buckets.

ii) Maintaining a fixed-size sample:

Suppose, we need to maintain a random sample S of size exactly s tuples.
e.g.: main memory size constraint.

Why? Do not know length of stream in advance.

Suppose at time n , we have seen n items.

- Each item is in the sample S with equal probability.

How to think about the problem: say $s=2$.

Stream: ~~already known~~ e.g.

At $n=5$, each of the first 5 tuples is included in the sample S with equal prob.

At $n=7$, each of the first 7 tuples is included in the sample S with equal prob.

Impractical solution would be to store all the n tuples seen. So, far and out of them pick s at random.

■ Solution: Fixed Size Sample

Reservoir Sampling Algorithm:

- Store all the first s elements of the stream S
- Suppose we have seen $n-1$ elements, and now the n^{th} element arrives ($n > s$).
 - With probability s/n , keep the n^{th} element, else discard it
 - If we picked the n^{th} element, then Good Write it replaces one of the s elements in the sample S , picked uniformly at random

Claim: This algorithm maintains a sample S with the desired property:

- After n elements, the sample contains each element seen so far with probability s/n .

Filtering of Data Stream: Filtering condition of a stream item is independent of other items of the same stream or any other data stream.

The most common example of such filtering is stream sampling, when each item is filtered out with a certain probability and the remaining items form the desired sample.

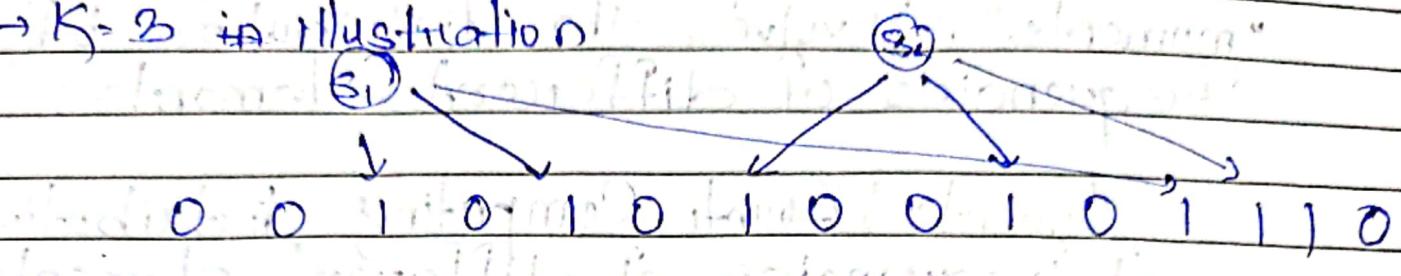
④ Bloom filter:

- i) A bloom filter consists of:
 - an array of n bits, initially all 0's.
 - A collection of hash function h_1, h_2, \dots, h_k .
- A set S of m key values.

- ii) The Bloom filter allows through all stream elements whose keys are in S , while rejecting most of the stream elements whose keys are not in S .

Illustration of Bloom Filters:

- Use K -independent hash functions instead of 1.
- $K=3$ in illustration



Counting Distinct Elements in a stream:

④ Flajolet-Martin Algorithm \Rightarrow

Flajolet-Martin Algorithm, also known as FM algorithm, is used to approximate the number of unique elements in a data stream off database in one pass.

Algorithm:

- Selecting a hash function h so each element in the set is mapped to a string to at least $\log_2 n$ bits.
- For each element x , $h(x)$ = length of trailing zeroes in $h(x)$.
- $R = \max(h(x))$
= Distinct elements = 2^R

Estimating Moments:

→ Estimating moments is a generalization of the problem of counting distinct elements in a stream. The problem, called computing "moments", involves the distribution of frequencies of different elements.

Moments: Counting distribution of frequencies of different elements in stream.

→ Suppose a stream consists of elements chosen from universal set. Assume the universal set is ordered so we can speak of the i^{th} element for any i .

→ Let, m_i be the number of occurrences of the i^{th} element for any i . Then the k^{th} -order moment of the stream is the sum over all i of $(m_i)^k$

Alon-Matias-Szegedy Algorithm: (AMS)

- AMS method works for all moments.
- Gives an unbiased estimate
- We will just concentrate on the 2nd moment of S .

Second moment: (Variance). It measures the spread of values of distribution or how far away from the normal.

→ Our goal is to compute $S = \sum_i n_i^2$

Algorithm:

- i) Choose a random hash function $h: \{1, 2, \dots, m\} \rightarrow \{1, 2, \dots, m^3\}$ from a pairwise independent hash family.
- ii) Initialize $z = 0$.
- iii) For each item x of the stream, update $z = \max(z_{\text{zeros}}(h(x)), z)$ where $z_{\text{zeros}}(y)$ denotes the number of trailing zeros of y in the binary representation.
- iv) Output 2^{z+e} for $e = 1/2$.

$$\boxed{E[f(x)] = 1/n \sum_{i=1}^n n(2e_i - 1)}$$

■ Counting oneness of a window?

→ Suppose we have a window of length N on a binary stream. We want at all times to be able to answer queries of the form "how many 1's are there in the last K bits". For any $K \leq N$.

→ To solve this problem, we use DGT algorithm.

→ DGT = Dasht - Grionis - Indyk - Motwani Algorithm.

→ DGIW Algorithm: Algorithm that uses $O(\log^2 N)$ bits to represent a window of N bits, and allows to estimate the number of 1's in the window with an error of no more than 50%.

- Each bit stream has a timestamp for the position at which it arrives.
- The first bit has timestamp 1, the second bit has 2 and so on.
- The positions are recognized with the window size N .
- The timestamp is represented with modulo N , and are represented as $\log_2 N$ bits.

The windows are divided into buckets consisting of -

- The timestamp at its right end
- the number of 1's must be in the power of 2 which are referred to as size of bucket.

Rules for forming Bucket:

- The right side of the bucket should always start with 1.
- Every bucket should have at least one 1, else no bucket can be formed.
- All bucket sizes should be a power of 2.

Updating the buckets:

- When a new bit comes in, chop the last bucket if its end-time is prior to N time units before the unit current time.
- If the current bit is 0, no other changes are needed.
- If the current bit is 1.
 - Create a new bucket of size 1. For just this bit, End timestamp = current time.
 - If there are now three buckets of size 1, Combine the oldest two into a bucket of size 2.

■ Decaying Window Algorithm: this algorithm allows you to identify the most popular elements in an incoming data streams

The decaying window algorithm not only tracks the most recutting elements in an incoming data stream, but also discounts any random spikes or spam requests that might have boosted an element's frequency.

- In decaying window, assign sum a score or weight to every element of the incoming data stream.
- Need to calculate the aggregate sum for each distinct element by adding all the weights assigned.

→ In decaying window algorithm, assign more weight to newer elements.

→ For new element, first reduce the weight of all the existing elements by a constant factor k .

→ Assign the new element with specific weight.

The aggregate sum of the decaying exponential weights can be calculated using the following formula:

$$\sum_{i=0}^{t-1} a t^{-i} (1-c)^i \quad | \quad c = \text{small constant of order } 10^{-6} \text{ to } 10^{-9}$$

- 1) Multiply the current sum/score by the value $(1-c)$.
- 2) Add the weight corresponding to the new element.

