

A.Y. 2024-25 (odd)

Department : CMBN

Semester : 7

Subject : Big Data Analytics (BDA)

Exam Date : 30/9/24

1A)

Columnar architecture style of Google Big Table :
Columnar architecture style of Google BigTable
allows it to:

- Store data in sparse multidimensional table
- Group related columns into column family, enabling efficient storage & retrieval
- Support high throughput for read/write operations, especially for large scale applications
- Effectively handle versions of data with timestamps.

This architecture makes BigTable ideal for scenarios like time-series data, web indexing, analytics & real-time monitoring.

(Any 2 points above can be elaborated in detail)

1B)

Comparison of Key-value pair architecture of NoSQL & with document oriented architecture of NoSQL.

P.T-O.

Space for
MarksQuestion
No.

START WRITING HERE

Aspect

Key-value

Document

1) Data model: Simple Key-value Documents
pair (JSON, XML, CSV)
with fields &
values

2) Query capabilities Look up by Key only Rich queries on fields, indexing & aggregation

3) Data Structure Unstructured (value can be anything) Semi-structured, hierarchical, supports nested data

4) Flexibility High flexibility with value formats but simple High flexibility with nested & complex documents

5) Use Cases Session data, Content management, simple lookups, Product catalogs etc. and user profiles.

Ans c) DMS architecture consists of following components.

- 1) Processor with standing queues
- 2) Limited working storage
- 3) Archival storage
- 4) I/O for stream data
- 4) Interface for ad-hoc query.

Space for
MarksQuestion
No.

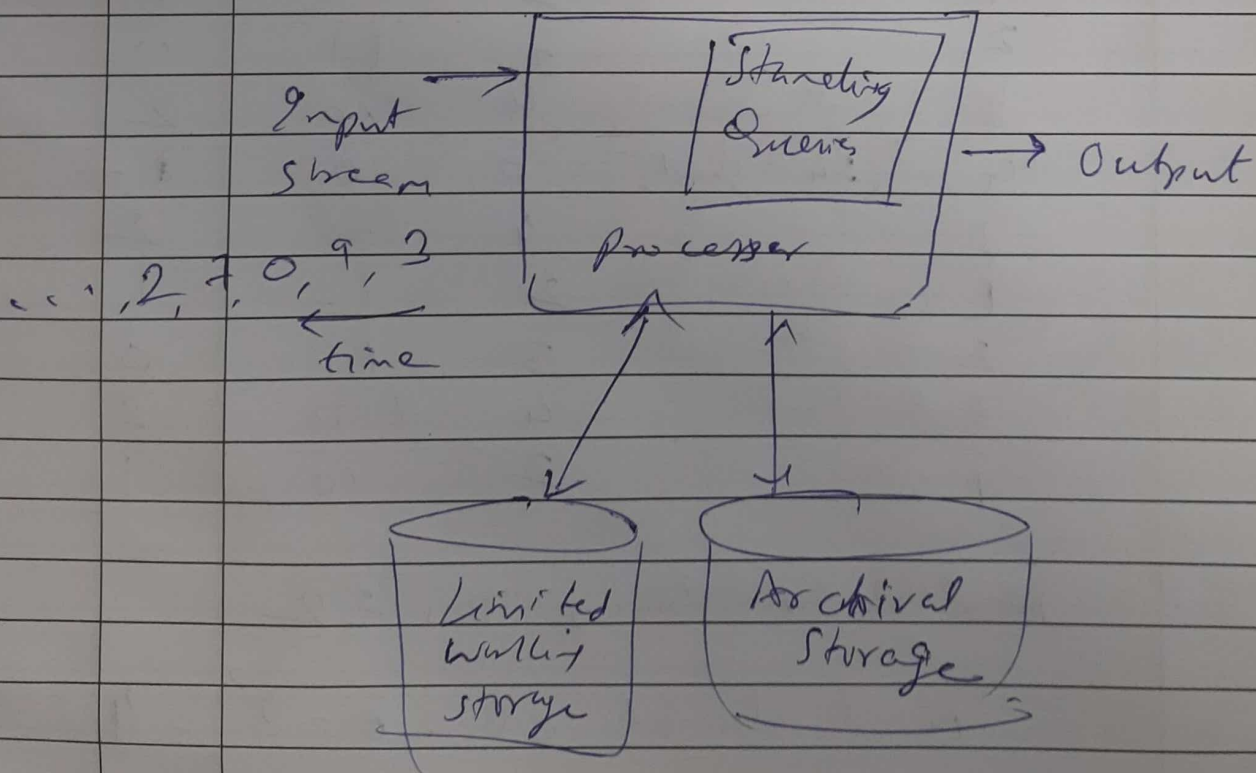
START WRITING HERE

Data sources can be IoT devices, sensors and other data integration entities. Inputs data can be in form of tuples like financial records etc..

Processor performs data cleaning, normalization & transformation to prepare data for processing.

Processed data required for standing queries are stored in limited working storage.

Large storage is stored in Archival storage that can be used to address Adhoc queries.



Space for
MarksQuestion
No.

START WRITING HERE

Standing queries are those that are asked about the data stream at all times.

Example: Report each new maximum value ever seen in stream S_1 .

Q2

(A)

No. of distinct elements in stream
 $S = 1, 3, 2, 1, 2, 3, 4, 3, 1, 2, 3, 1$
 $h(x) = (6x + 1) \bmod 5$

Date	Binary Hash	Binary	Trailing 0's
1	$(6+1) \bmod 5 = 2$	10	1
3	$(18+1) \bmod 5 = 4$	100	2
2	$(12+1) \bmod 5 = 3$	11	0
1	2	10	1
2	3	11	0
3	4	100	2
4	$(24+1) \bmod 5 = 0$	0	1 →
3	4	100	2
1	2	10	1
2	3	11	0
3	4	100	2
1	2	10	1

Max no. of trailing 0's = 2

∴ No. of distinct elements = $2^2 = 4$

Space for
MarksQuestion
No.

START WRITING HERE

2 B) Bloom Filter :-

Given a set list of elements (S),
we need to determine whether an
element x is in S or not?

Solution using Bloom's Filter

Consider $|S| = m$ & $|B| = n$
where B is hash table

Initially set all bits of B to 0s

For each element ~~y~~ $y \in S$ use
hash function h_i for $i = 1, \dots, k$
and set $B[h_i(y)] = 1$. for
each hash function

When an element x arrives

If $B[h_i(x)] = 1$ for all
 $i = 1, 2, \dots, k$

then declare that x is in S .

(i.e. if x hashes to all locations of
~~bit 1~~ that are 1,
then most likely x was seen
earlier).

Otherwise discard x , as not in S .

Space for
MarksQuestion
No.

START WRITING HERE

Example Consider set of elements as
3, 5, 2, 1

hash functions

$$h_1 = (3x + 1) \bmod 11$$

$$h_2 = (6x + 3) \bmod 11$$

$$h_3 = (5x + 2) \bmod 11$$

Bucket size = 11

element 3 hashes to

$$h_1(3) = 10 \bmod 11 = 10$$

$$h_2(3) = 21 \bmod 11 = 10$$

$$h_3(3) = 17 \bmod 11 = 6$$

0	1	2	3	4	5	6	7	8	9	10
						1				1

Similarly, other elements

	$h_1(k)$	$h_2(k)$	$h_3(k)$
5	5	0	5
2	7	4	1
1	4	9	7

Final Bucket

0	1	2	3	4	5	6	7	8	9	10
1	1	0	0	1	1	1	1	0	1	1

Consider element 4.

$$h_1(4) = 13 \bmod 11 = 2$$

\therefore Slot 2 is 0 \therefore 4 is not seen in set.

Space for
MarksQuestion
No.

START WRITING HERE

2c) The most suitable technique for the E-commerce website is to use DGM algorithm.

The DGM (Data-Gionis-Indyk-Motwani) algorithm is an efficient method for tracking the number of 1's in the last N bits of a binary stream using limited memory.

The E-commerce website can maintain a separate bit stream for each product. A transaction in which a product is present is set to 1 and if absent then the bit is assumed as 0.

For each '0' in the stream, DGM does nothing. It creates a new window of size L only when the stream has 1. Hence

Hence, considering the "sparsity of data", DGM would be suitable as it will create a new window of size L only for the set of products involved in the transaction.

The algorithm also merges smaller buckets & hence uses less complexity $O(\log^2 N)$ to track & count 1's in last N bits.

Space for
MarksQuestion
No.

START WRITING HERE

3A)

The most suitable filtering technique is Content based ~~filter~~ recommendation system. It can be used to suggest songs to customers by analyzing the characteristics of content of the songs themselves, rather than relying on other's preference as in collaborative filtering. It focuses on matching songs with similar attributes to those the user has already liked.

The first step for content based recommendation for Spotify would be to define and extract features. Each song is broken down into features such as genre of the song, singer/artist, music instruments used, theme of the song, melody, rhythm, tempo, language, etc.

Step 2 is to create a user profile based on the songs that user has rated or heard. The profile contains weighted list of features for these songs.

Ex: User 1 [Romantic, Acoustic, Sad, Kind, ...] having weights like $\langle 0.4, 0.8, 0.7, 0.9, \dots \rangle$.

Space for
MarksQuestion
No.

START WRITING HERE

Step 3 is to compare new songs with user's profile i.e. finding songs having similar features. This can be achieved by using cosine or similarity.

Final step is to rank songs based on their similarity score & suggest songs with high similarity with user's profile.

This technique will have advantage of personalization & can overcome the problem of cold start i.e. recommend the user song of a singer Arjit Singh, which is having features like Romantic & Sad & in Hindi language to users whose profile matches with the new released song.

Limitations.

- Over specialisation: User may only get recommendations based on their previous ratings thereby not giving exposure to user of other genre or artists.
- Feature Engineering: It relies on system's ability to extract meaningful features of the songs.

Space for
MarksQuestion
No.

START WRITING HERE

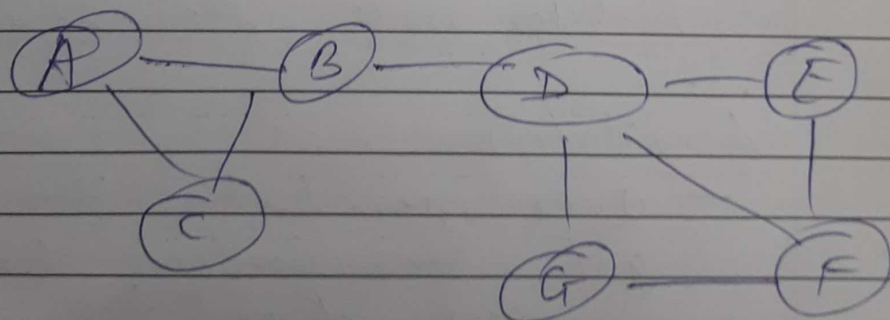
3B)

A most suitable technique for this social media platform to detect communities is to use Girvan-Newman (GN) algorithm.

GN algorithm progressively removes edges that are most likely acting as a bridge between different communities.

To determine the edge that connects 2 communities, GN algorithm ~~computes~~ edge betweenness. The edge that has highest edge betweenness score is considered to be the bridge of 2 different communities and hence removed/discarded.

Consider a social graph below:



Step 1: Apply modified BFS. In this BFS, a node of level i that is connected to more than

Space for
Marks

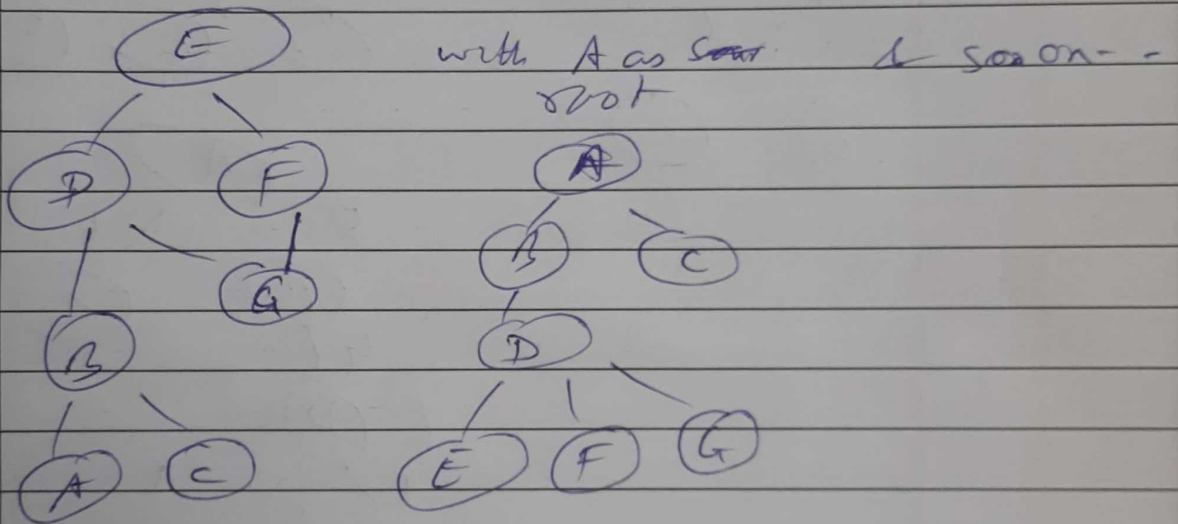
Question
No.

START WRITING HERE

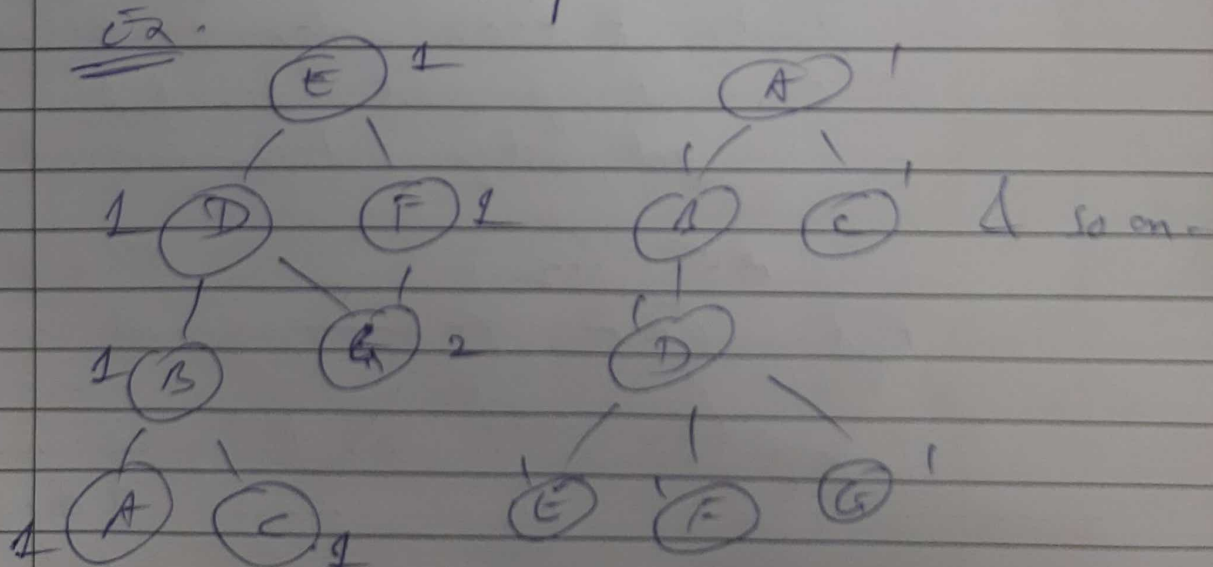
one node of level i , is considered
to be child / reachable from
both nodes.

Hence output of this modified BFS
is not a spanning tree but a DAG.

For the given example, with E as the
root node, we get.



~~1st~~ Step 2: It labels the nodes with
number of paths from the root.

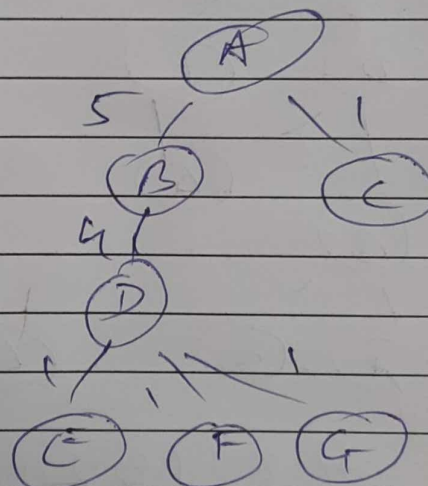
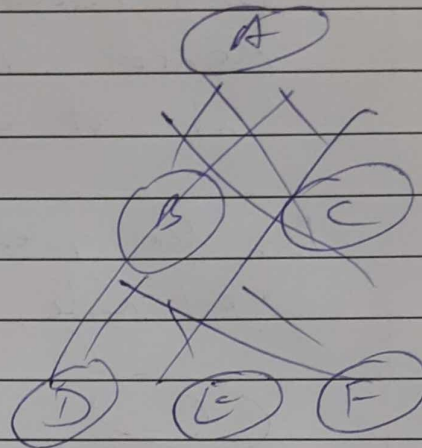
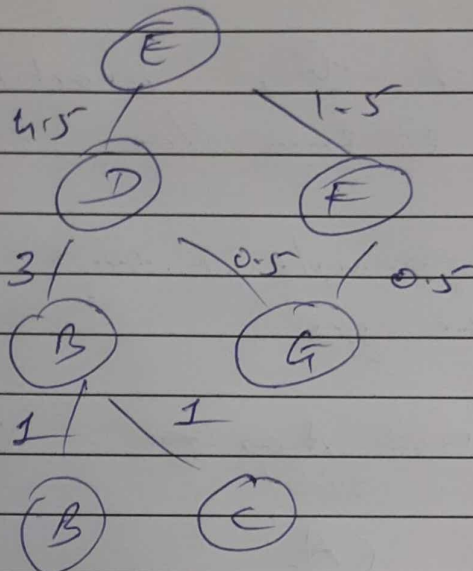


Space for
Marks

Question
No.

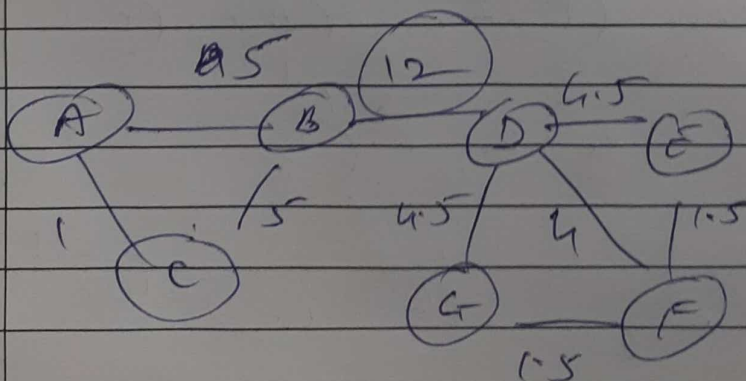
START WRITING HERE

~~First~~ Step 3: It calculates credit from
leaf to root. For example.



4 score

The above steps after repeating for all
roots, their edge contribution is
summed up and divided by 2
to get edge between nodes. For example:



The edge BD
is discarded
to get 2
communities
{A, B, C} &
{D, E, F, G}.