USN | | | | | | | 0 | 4 | 3 |

18CS72

## Seventh Semester B.E. Degree Examination, Dec.2023/Jan.2024
### Big Data and Analytics

Time: 3 hrs.

Max. Marks: 100

**Note: Answer any FIVE full questions, choosing ONE full question from each module.**

### Module-1

1. a. How is Data Architecture layers used for analytics? Explain with functions of each layer. **(10 Marks)**
   b. Briefly describe the three fundamental services offered by Cloud Computing. **(10 Marks)**

**OR**

2. a. List the features of Grid Computing. How does it differ from clusters and cloud computing. **(10 Marks)**
   b. Why is Data quality important in discovering new knowledge and decision making? **(10 Marks)**

### Module-2

3. a. List Hadoop core components and explain with appropriate diagram. **(10 Marks)**
   b. Explain the working of the Hadoop Map Reduce framework. **(10 Marks)**

**OR**

4. a. Explain the working of Hadoop – 2 Execution model (YARN Model). **(10 Marks)**
   b. With a diagram, explain the concept of APACHE Sqoop to acquire relational data. **(10 Marks)**

### Module-3

5. a. Define NOSQL Explain Big Data NOSQL or Not – only SQL with its features, transactions and solutions. **(10 Marks)**
   b. Describe graph database characteristic, typical used and examples. **(10 Marks)**

**OR**

6. a. Explain Mongo DB with its features. **(10 Marks)**
   b. Compare and contrast RDBMS and Mongo DB databases. **(05 Marks)**
   c. What are the different ways of handling Big Data Problems? **(05 Marks)**

### Module-4

7. a. Describe the Hive architecture components along with Hive Built – in functions. **(10 Marks)**
   b. Explain with respect to Hive QL
   i) Hive QL Data Definition Language (DDL).
   ii) Hive QL Data Manipulation Language (DML). **(10 Marks)**

**OR**

8. a. Explain the architecture, feature and applications of PIG. **(10 Marks)**
   b. Illustrate by considering an example the working of the Map Reduce programming model. **(10 Marks)**

## Module-5

9  a. How does regression analysis predict the value of the dependent variable in case of linear regression? **(10 Marks)**

b. Explain with example and algorithm, the working principle of Apriori process for adopting the subset of frequent item sets as a frequent itemset. **(10 Marks)**

**OR**

10  a. Define Web Mining. Discuss the broad classification of web mining and their applications. **(10 Marks)**

b. Define the term Social network. Explain social network as graphs with Centralities, Ranking and Anomaly Detection. **(10 Marks)**

| | **VTU Examination – Dec- 23/Jan- 24** <br> **Solution** | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Sub: | **Big Data and Analytics** | | | | Sub Code: | **18CS72** | Branch: | **ISE** |
| Exam Date: | **31/01/2024** | Duration: | **3 Hrs** | Max Marks: | **100** | Sem | **VII** | |

| | **Answer any FIVE FULL Questions** | **MARKS** | **CO** |
|---|---|---|---|
| | **MODULE-1** | | |
| 1 (a) | a. How is Data Architecture layers used for analytics? Explain with functions of each layer <br> **Solution:** <br><br>  <br><br> Fig. w.r.t. logical layers in a data architecture ——— 5M <br> Function in detail ——— 5M | [10] | CO1 |
| (b) | **Briefly describe the three fundamental services offered by Cloud Computing.** <br> **Solution:** <br> Cloud services can be classified into three fundamental types i) Infrastructure as a Service(IaaS) ii) Platform as a Service(PaaS) iii) Software as a Service(SaaS) <br><br> Cloud computing –Definition <br> ——— features –4M | [10] | CO1 |

*Fundamental types — IaaS — 2 M*
*to be explained* } — PaaS — 2 M
*in detail including* — SaaS — 2 M
*Examples*

**OR**

2 (a) **List the features of Grid Computing. How does it differ from clusters and cloud computing.** [10] CO1
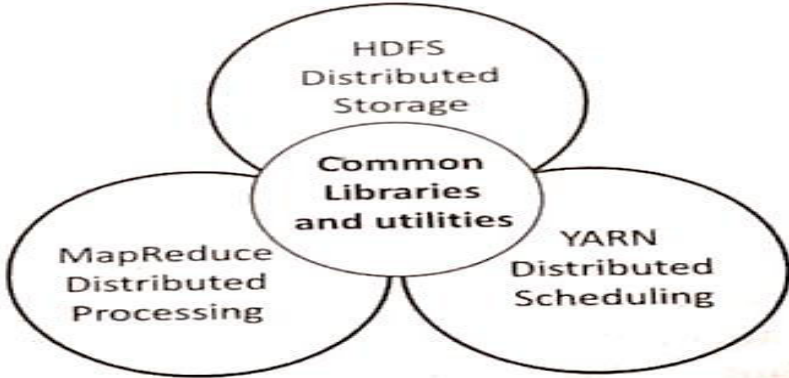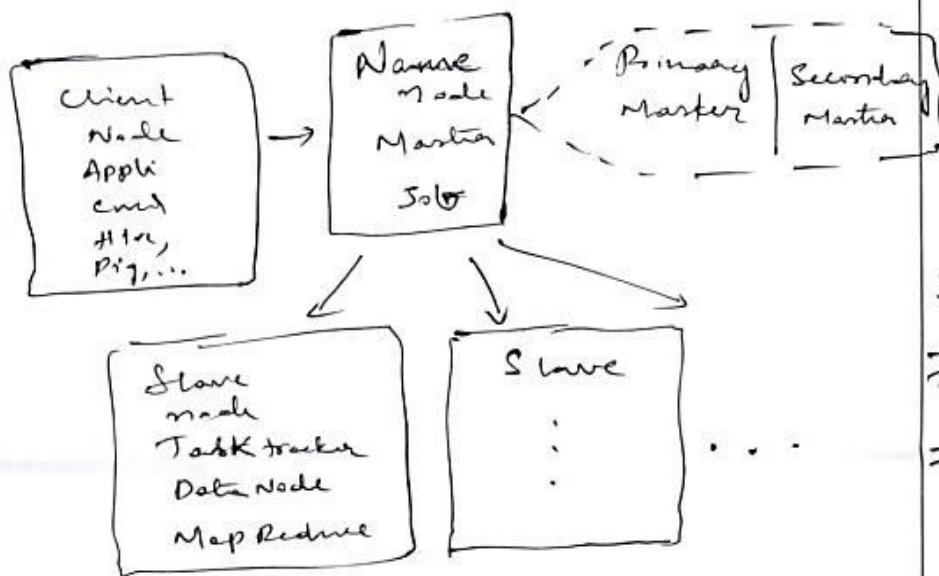
**Solution:**

Features of Grid Computing
• Grid computing is similar to cloud computing ,it is scalable.
• Cloud computing depends on sharing of resources (for eg , networks , servers ,storage , applications and services )to attain coordination and coherence among resources similar to grid computing.
• Grid also forms a distributed network for resource integration.     **4 Marks**

| | Cluster Computing | Grid Computing | Cloud Computing |
|---|---|---|---|
| **Basic Idea** | Aggregation of resources. | Segregation of Resources. | Consolidation of Resources. |
| **Running Processes** | Same processes run on all computers over the cluster at the same time. | Job is divided into sub-jobs each is assigned to an idle CPU so they all run concurrently. | Depends on service provisioning. Which computer offers a service and provisions it to the requesting clients. |
| **Operating System** | All nodes must run the same operating system. | No restriction is made on the operating system. | No restriction is made on the operating system. |
| **Job Execution** | Execution depends on job scheduling. So, jobs wait unit it's assigned a runtime. | Execution is scalable in a way that moves the execution of a job to an idle processor (node). | Self-Managed. |
| **Suitable for Apps** | Cascading tasks. If one tasks depends on another one. | Not suitable for cascading tasks. | On-demand service provisioning. |
| **Location of nodes** | Physically in the same location | Distributed geographically all over the globe. | Location doesn't matter |
| **Homo/Heterogeneity** | Homogenous | Heterogeneous | Heterogeneous |
| **Virtualization** | None | None | Virtualization is a key |
| **Transparency** | Yes | Yes | Yes |
| **Security** | High | High, but doesn't reach the level of cluster computing. | Lower than both types. |
| **Interoperability** | Yes | Yes | No |
| **Application Domains** | industrial sector, research centers, health care, and centers that offer services on the nation-wide level | industrial sector, research centers, health care, and centers that offer services on the nation-wide level | Banking, Insurance, Weather Forecasting, Space Exploration, Business, IaaS, PaaS, SaaS |
| **Implementation** | Easy | Difficult | Difficult – need to be done by the host. |
| **Management** | Easy | Difficult | Difficult |
| **Resource Management** | Centralized (locally) | Distributed | Both centralized and distributed. |
| **Internet** | No internet access is required | Required | Required |

                                                                                                    **6 Marks**

| | | | |
|---|---|---|---|
| (b) | Why is Data quality important in discovering new knowledge and decision making?<br>Solution: | [10] | CO1 |

Definition of Data quality

— 2M

Five R's    Relevancy
Recency
Range
Robustness
Reliability

Removing Data Noise, Outliers,
Missing & Duplicate values. with
detailed explanation & Justification
for Knowledge & Decision making

— 8M

### MODULE-2

| 3 (a) | List Hadoop Core Components and explain with appropriate diagram.<br>Solution: | [10] | CO2 |



**Diagram + Explanation – 2+8 marks**

| (b) | Explain the working of the Hadoop MapReduce Framework.<br>Solution: | [10] | CO2 |

The client, Master, Namenode &
Slave nodes with detailed explanation
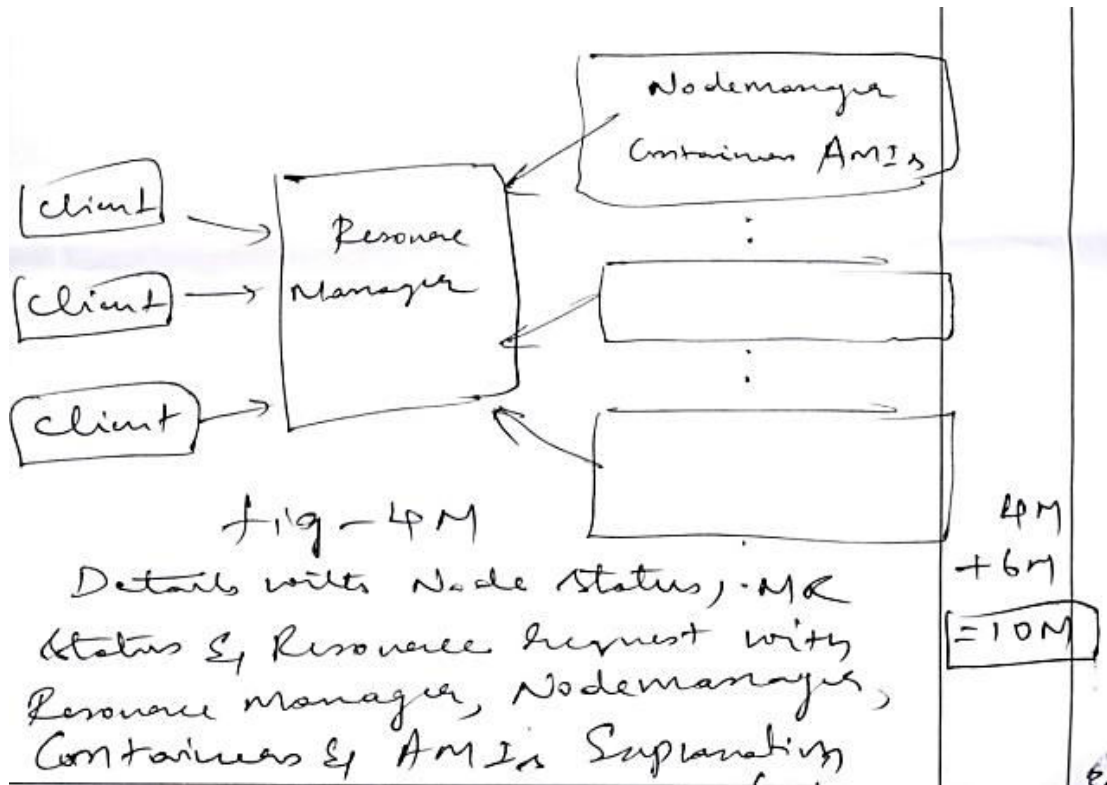fig — 3M
Explanation — 7M

3M +
7M
= 10M

**OR**

| 4 (a) | **Explain the working of Hadoop-2 Execution model(YARN model)**<br>**Solution:** | [10] | CO2 |
|---|---|---|---|



fig – 4M

Details with Node status, MR
status & Resource request with
Resource manager, Nodemanager,
Containers & AMIs Explanation

4M
+ 6M
= 10M

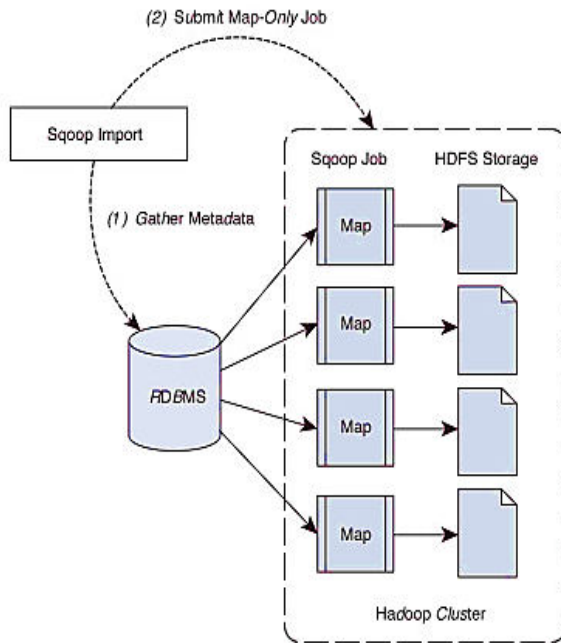| (b) | **With a diagram, explain the concept of APACHE sqoop to acquire relational data.**<br>**Solution:**<br><br>Diagram – 4 Marks<br>Explanation – 6 Marks | [10] | CO2 |
|---|---|---|---|

Fig: Two-step Apache Sqoop data import method.          Fig: Two-step Sqoop data export method.

## MODULE-3

| | | |
|---|---|---|
| 5 (a) **Define NOSQL. Explain Big Data NoSQL or Not-only SQL with its features, transactions and solutions.**<br>**Solution:**<br><br>Definition — 2M<br><br>Features: Relax of ACID Properties,<br>Two properties of CAP theorem<br>BASE thrive model — 4M<br><br>Solutions: Apache's HBase, MongoDB,<br>Cassandra, CouchDB, Oracle & Riak with<br>uses.  — 4M | [10] | CO3 |
| (b) **Describe graph database characteristic, typical uses and examples.**<br>**Solution:**<br><br>1 Characteristics: Specialized query lang,<br>(XDF), different models, hyper-edges,<br>Joins, ...  — 4M<br><br>uses: link analysis, friend of friend<br>queries, Rules & Inference, Rule Induction,<br>pattern matching  — 4M<br><br>Examples: Neo4J, Allegro graph, Hyper graph,<br>Infinite, Titan & Flock DB.  — 2M | [10] | CO3 |

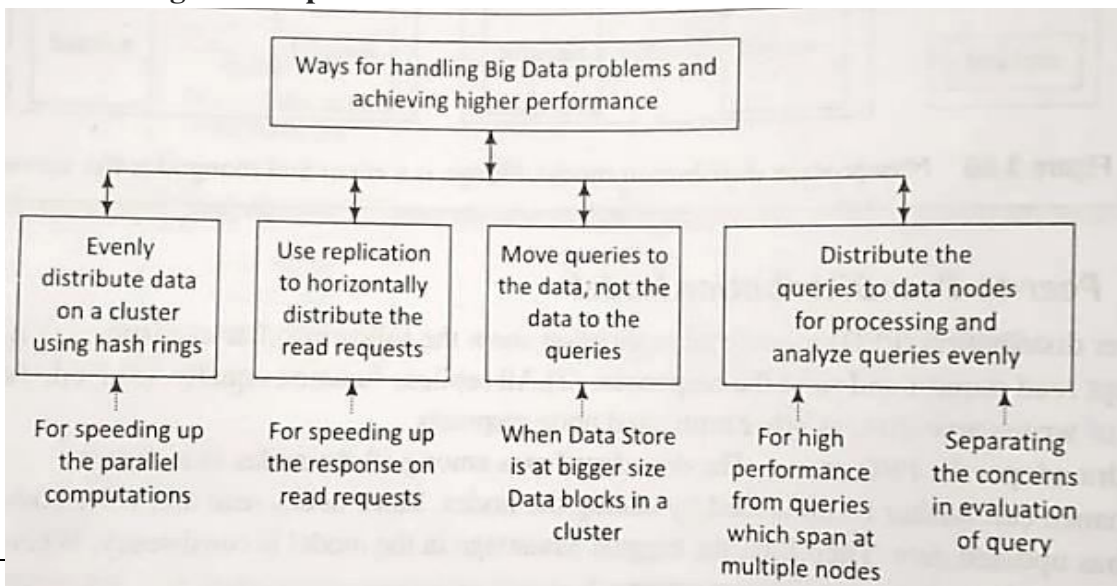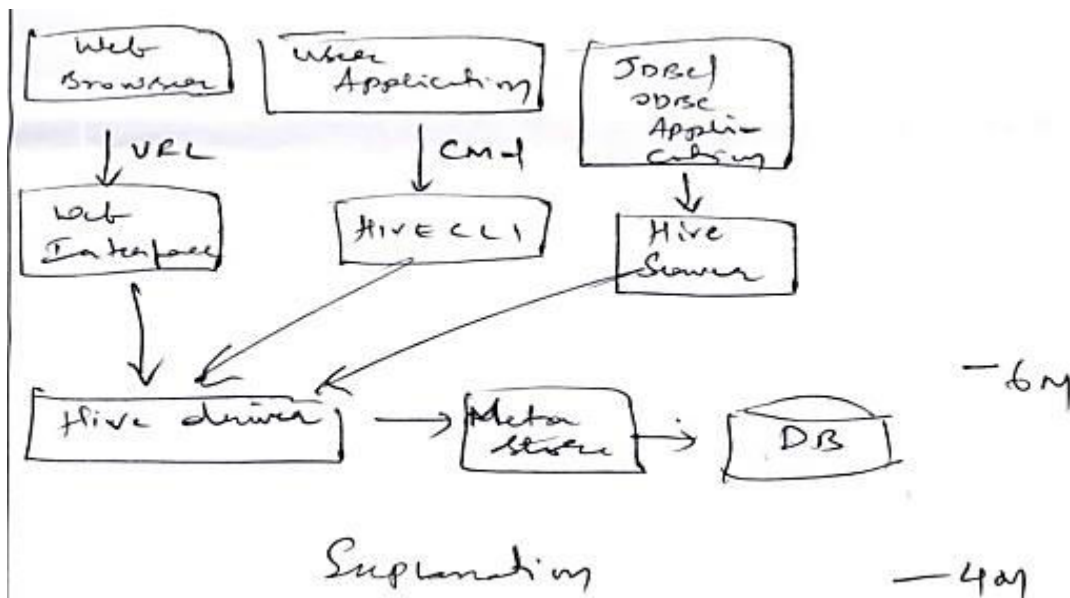|  |  | **OR** |  |  |
|---|---|---|---|---|
| 6 (a) | **Explain MongoDB with the features.**<br>**Solution:**<br><br>MongoDB is non-relational, NoSQL, Distributed, open source, document based, Cross platform, Scalable, flexible, Indexed, Multimaster / fault tolerant.<br>Features: Physical Container for Collections<br><br>Collection stores, well defined, JSON style documents, BSON Serialization format, Efficient, Deep query ability, No complex Joins, Distributed DB, Indexed, Atomic operations, Fast in place updates, No Configurable cache, conversion/mapping of application objects, ...<br><br>4M + 6M = 10M | [10] | CO3 |
| (b) | **Compare and Contrast RDBMS and Mongo DB databases.**<br>**Solution:**<br><br>Comparing with Features like Data Model, Schema, Typed Data, locality, updates, Transactions, Auditing & Scaling<br><br>5M | [05] | CO3 |
| (c) | **What are the different ways of handling Big Data Problems?**<br>**Solution: Diagram+ Explanation – 2+3 Marks** | [05] | CO3 |

| | | | | |
|---|---|---|---|---|

**7 (a)** **Describe the Hive architecture components along with Hive Built-in functions.** [10] CO4

**Solution:**



**(b)** **Explain with respect to HiveQL.** [10] CO4
   i) **Hive QL Data Definition Language(DDL)**
   ii) **Hive QL Data Manipulation Language (DML).**

**Solution:**

i) a) CREATE DB.
   b) SHOW DB
   c) CREATE Schema
   d) CREATE TABLE
   } HIVE QL DDL Commands 5M

ii) a) DROP DB
   b) DROP Schema
   c) ALTER TABLE
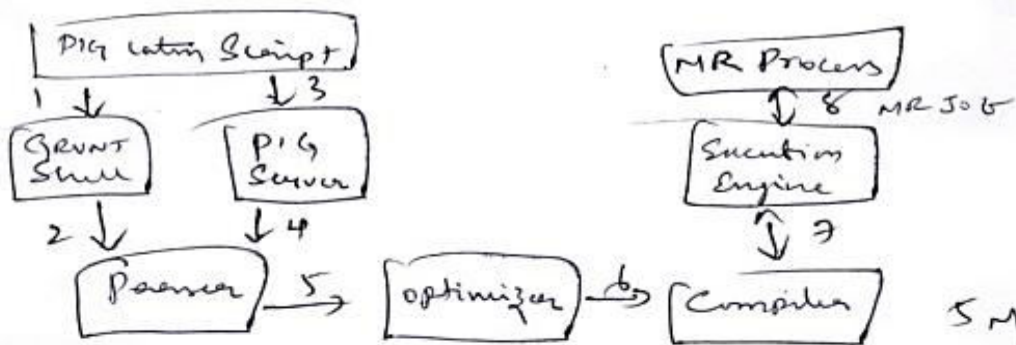   d) DROP table
   e) LOAD Data
   } Hive QL DML 5M Commands

**OR**

| | | | |
|---|---|---|---|
| 8 (a) | **Explain the architecture, feature and applications of PIG.**<br>**Solution:** | [10] | CO4 |



Architecture — 5M

Features : Scripts, UDFs, Variety of data, multi query approach, Inconsistent Schema, ETL, auto optimizations, deals with whole operations, Fewer codes, live anywhere, take anything, run as if flying. — 3M

applications : Analysing large datasets, adhoc processing, large data sources, supports different data types, processing time sensitive data loads. — 2M

| | | | |
|---|---|---|---|
| (b) | **Illustrate by considering an example the working of the Map Reduce Programming model.**<br>**Solution:** | [10] | CO4 |



— 4M

Diagram with detailed explanation
— 6M

| | | |
|---|---|---|
| 8 (a) | **MODULE-5** | |

| | | | | |
|---|---|---|---|---|
| 9 (a) | **How does regression analysis predict the value of the dependent variable in case of linear regression?**<br>**Solution:** | [10] | CO5 |
| | | Simple linear Regression with<br><br>Equations & Explanations    10M | | |
| (b) | **Explain with example and algorithm, the working principle of Apriori process for adopting the subset of frequent item sets as a frequent Itemset.**<br>**Solution:** | [10] | CO5 |

Algorithm

$C_k$: Set of Candidate -Itemset

$F_i$: Set of frequent itemset

$F_1$ = {Large items}

for c     ) do {

$C_{k+1}$ = New candidates from $F_k$

for each transaction t do

$C_{k+1}$    } New calculations
$F_{k+1}$

                   — 3M

Example:                 Explanation — 2M

| TID | Items |
|---|---|
| 1 | {A, C, D} |
| 2 | {A, B, C, E} |
| 3 | {B, E} |
| 4 | {B, C, E} |

⇒ Several Steps       — 3M

3M<br>+2M<br>+3M<br>+2M<br><br>=10M

Real result    ⇒

| Itemset | Support |
|---|---|
| {B, C, E} | 2 |

any Equivalent Example with Explanation —2M

| | | **OR** | | |
|---|---|---|---|---|

| 10 (a) | **Define web mining. Discuss the broad classification of web mining and their applications.**<br>**Solution:** | [10] | CO5 |

Web mining

Web Content mining      Web Structure mining      Web usage mining

Text<br>Image<br>Video<br>audio<br>:

Hyper links    Document storage

* Intra<br>* Inter

* Web server logs    10M<br>* Appli. —<br>* Appli. server logs
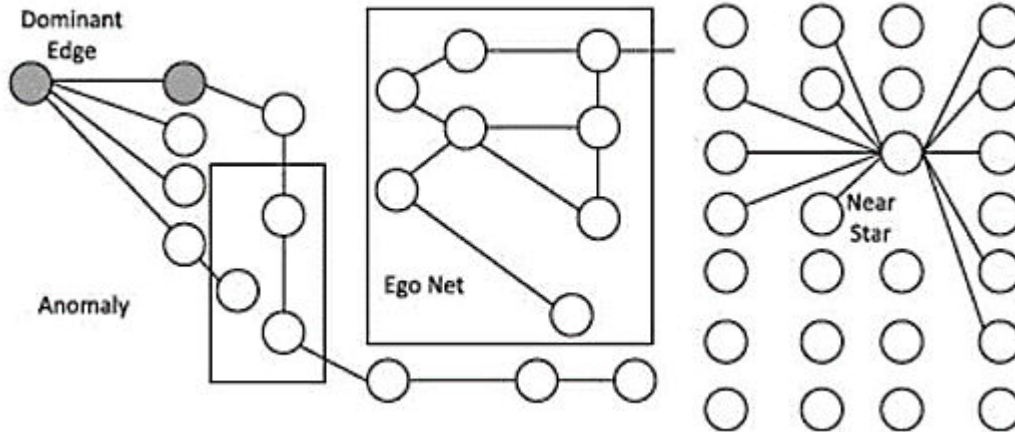
Proper Classifications & applications

(b) **Define the term Social network. Explain social network as graphs with centralities, Ranking and Anomaly Detection.**

**Solution:**

A social network is a social structure made of individuals (or organizations) called "nodes," which are tied (connected) by one or more specific types of inter-dependency, such as friendship, kinship, financial exchange, dislike or relationships of beliefs, knowledge or prestige.



Definition — 2M
Centralities — 2M
Ranking — 2M
Anomaly Detection — 4M

[10]  CO5