

(2 Hours.)

[Total Marks: 60]

- N. B.:** (1) All questions are compulsory.  
 (2) Numbers to the right indicate marks.  
 (3) Make suitable assumptions wherever necessary and state the assumptions made.  
 (4) Answers to the same question must be written together.  
 (5) Mixing of Sub-Questions is not allowed.  
 (6) Draw neat labelled diagrams wherever necessary.

**1. Attempt any two of the following:**

**12**

- What is Big Data? State its characteristics.
- Discuss the challenges of Big data.
- Describe the current analytical architecture for data scientists.
- What are different phases of the Data Analytics Lifecycle? Explain in detail.

**2. Attempt any two of the following:**

**12**

- What is clustering? State its advantages & disadvantages.
- Explain the concept of K-means clustering algorithm.
- Write a brief note on Apriori algorithm.
- What is Logistic regression? Explain in detail.

**3. Attempt any two of the following:**

**12**

- Write a short note on Decision trees.
- Explain a probabilistic classification method based on Naive Bayes' theorem.
- Explain Time Series Analysis in detail.
- What is tokenization? Explain how it is used in text analysis.

**4. Attempt any two of the following:**

**12**

- Describe the concept of MapReduce.
- State the components of Spark Stack.
- Write a short note on Hadoop Distributed File System.
- Explain in detail the Hadoop Architecture.

**5. Attempt any two of the following:**

**12**

- What is the functionality of the filter mapper? Explain in detail with example.
- Write a short note on Data Ingestion.
- Explain with example Relations, tuples and Filtering in context of Pig.
- Write the entire procedure with appropriate commands for importing data from MySQL to HDFS.

\*\*\*\*\*