

M. Sc. DEGREE EXAMINATION, APRIL 2022
BIOINFORMATICS
FOURTH SEMESTER

COURSE : CORE
PAPER : BIG DATA ANALYSIS
TIME : 3 HOURS

MAX. MARKS: 100

SECTION – A

ANSWER ALL QUESTIONS

(20 X 1=20)

Choose the correct answer:

- Which is the correct order for medical image analysis
 - Preprocessing, Segmentation, Thresholding, and feature detection
 - Preprocessing, Thresholding, Segmentation, and feature detection
 - Preprocessing, Segmentation, feature detection and Thresholding
 - Preprocessing, feature detection, Segmentation and Thresholding
- Handling large volumes of data involves
 - Using correct algorithms
 - Choosing the right data structures
 - Using right tools
 - All the above
- Volume, velocity and variety are ____ to big data,
 - Intrinsic
 - Extrinsic
 - Both a and b
 - none of the above
- Which is not an discretization technique?
 - Normalization
 - Histogram
 - Binning
 - Regression
- _____ function is responsible for consolidating the results produced by each of the Map() functions/tasks.
 - Reduce
 - Map
 - Task
 - All of them

Fill in the Blanks:

- _____ is the process of reducing a group of words into their lemma or dictionary form in NLP.
- _____ is the python library will take care of the matrix splitting, and linear regression variable weights can be calculated using matrix calculus.
- Variety describes one of the biggest challenges of _____.
- PCA stands for _____.
- ____ and ____ are the two ways to implement DFS.

Define in single line:

- Define IoT.
- Give an example of graph database
- What is the five step process to structure your data analysis?
- Define Binning.
- List two features of DFS.

Say True or False:

- Dummy variables* can only take two values: true(1) or false(0).
- Is sparse matrix also known as dense matrix?
- Veracity makes sure that the data is accurate?
- Is smoothing technique can be employed to remove the noise in data?
- Blocks, name node and data node are the concepts of HDFS?

SECTION – B

ANSWER ANY FOUR QUESTIONS. EACH ANSWER SHOULD NOT EXCEED 500 WORDS. ALL QUESTIONS CARRY EQUAL MARKS. (4 x 10 = 40)

21. Comment on implementation of Network modelling in biological concepts.
22. Elaborate the data integration process briefly.
23. Explain the 3V's of big data.
24. Highlight the importance of data reduction process with PCA as example.
25. Brief about the following: a) Data Scaling, b) Data Transformation.
26. Discuss the key points in analyzing the data and communicating the results.
27. Illustrate the importance of MapReduce.

SECTION – C

ANSWER ANY TWO QUESTIONS. EACH ANSWER SHOULD NOT EXCEED 1200 WORDS. ALL QUESTIONS CARRY EQUAL MARKS. (2 x 20 = 40)

28. Illustrate the facets of big data in detail.
29. Discuss the machine generated data and their advantages
30. Enumerate the data per-processing steps in detail.
31. Comment on HDFS and YARN as resource manager.
