

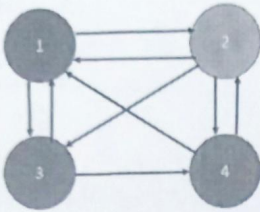
DECEMBER 2021: END SEMESTER ASSESSMENT (ESA) B TECH V SEMESTER

UE19CS322 – BIG DATA

Time: 3 Hrs	Answer All Questions Clearly mark the question number/part and write the entire answer for a subquestion together.	Max Marks: 100
-------------	---	----------------

1	a)	Volume, Veracity, Variety and Velocity are considered to be the 4V's of Big Data. Give an example to demonstrate each of these. Also, give two reasons as to why Volume has become a challenge in Big Data when compared to traditional data.	6
	b)	"Combiners are considered as mini-reducers in Map Reduce" – What are combiners and why are they considered as mini reducers. A map reduce programmer submitted a MR program that had optimized the MR program by replacing the reducer by a combiner. Is this guaranteed to give correct results? Justify	4
	c)	What is HDFS – briefly explain any two features of HDFS? Why is the HDFS architecture termed a master-slave architecture – briefly explain the major different components and their functionality? Discuss the motivation behind this architectural organization? (3+5+2)	10
2	a)	1. Which of the following components does not belong to the Hadoop eco-system? a. Sqoop b. Pig c. Mahout d. Dbase - 2. Workflow does not involve a. Flow of work from consumer to producer b. Steps to be run c. Executing the specified steps identified - d. Error handling 3. Apache Oozie a) Schedules the different types of jobs b) Runs different types of jobs c) Manages different types of jobs - d) All the above 4. Identify the odd one. These are needed to install Ambari in a Cluster a) Stack of Services b) Services - c) Components making up the service d) CLI for installation of the clusters	4

- b) Write down the adjacency matrix and transition matrix (columns as source and rows as destination) of the following directed graph and compute page rank using Map Reduce for 1 iteration assuming that initial page ranks of each node is $\frac{1}{4}$. Please use the format shown below to show the working of the Map reduce program. Assume that the transition matrix is stored on HDFS in two blocks – columns 1 and 2 are stored on Node 1 and columns 3 and 4 on Node 2



Node 1 Map Intermediate key	Node 1 Map Intermediate value	Node 2 Map Intermediate key	Node 2 Map intermediate values
Reduce input key	Reduce input value	Reduce output key	Reduce output values

- c) Given that you have the following data stored on HDFS using CSV files
TaxPaid Table

Name	PAN #	Date	Tax Paid

BankDetails

PAN #	State in which tax paid	Bank account.

You need to design a Map-reduce program to compute the total tax paid per state. Which relational operations will you require to use to perform this computation. How many map-reduce steps will you require to perform the computation (2 marks) and what will be relational operation performed? (4 marks) Express your solution in the following format

MR Step no	Relational Operational operations performed	Intermediate key value pairs; Intermediate keys are output of mapper	Final key-values; final keys are output of reducer.

