

Name:	 UPES UNIVERSITY WITH A PURPOSE
Enrolment No:	

UNIVERSITY OF PETROLEUM AND ENERGY STUDIES
End Semester Examination, May 2021

Course: In Memory Processing
Program: B.Tech CSE Big Data
Course Code: CSBD 3003

Semester: VI
Time 03 hrs.
Max. Marks: 100

SECTION A

- 1. Each Question will carry 5 Marks**
2. Instruction: Write brief answers

S. No.		CO
Q 1	Discuss the role of Catalyst Optimizer in Spark SQL?	CO3
Q 2	Explain the concept of Resilient Distributed Dataset (RDD) and how do we create RDDs in Spark?	CO2
Q 3	Differentiate between CreateOrReplaceTempView and createGlobalTempView?	CO1
Q 4	What do you understand by Pair RDD?	CO2
Q 5	What are the components of Spark Ecosystem?	CO1
Q 6	Distinguish between map and flatMap transformation in Spark?	CO2

SECTION B

- 1. Each question will carry 10 marks**
2. Instruction: Write short notes

Q 7	Define shuffling in Spark. When does it occur ? Does shuffling change the number of partitions?	CO2
Q 8	Distinguish between wide transformation and narrow transformation with suitable examples.	CO1
Q 9	Does Spark SQL help in big data analytics through external tools too? Justify	CO2
Q 10	When do you use apache spark and what are the benefits of Spark over Mapreduce framework?	CO1
Q 11	What are the different storage levels of persistence in Spark?	CO3

SECTION-C

1. Each question carries 20 Marks.

2. Instruction: Write long answer.

Q 12	<p>IMBD is an online database of movie-related information. IMBD users rate the movies and provide reviews. They rate the movies on a scale of 1 to 5; 1 being the worst and 5 being the best. The dataset also has additional information, such as the release year of the movie. You have to analyse the data collected and answer the following questions.</p> <p>Sample data from the dataset: (Title, release year, rating, number of users)</p> <p>1,The Nightmare Before Christmas,1993,3.9,4568 2,The Mummy,1932,3.5,4388 3,Orphans of the Storm,1921,3.2,9062</p> <p>You need to find:</p> <table style="width: 100%;"><tr><td>1) The total number of movies</td><td style="text-align: right;">(2)</td></tr><tr><td>2) The maximum rating of movies</td><td style="text-align: right;">(2)</td></tr><tr><td>3) The number of movies that have maximum rating</td><td style="text-align: right;">(3)</td></tr><tr><td>4) The movies with ratings 1 and 2</td><td style="text-align: right;">(3)</td></tr><tr><td>5) The list of years and number of movies released each year</td><td style="text-align: right;">(5)</td></tr><tr><td>6) The number of movies that have a runtime of two hours</td><td style="text-align: right;">(5)</td></tr></table> <p style="text-align: center;">OR</p> <p>a) Discuss Broadcast variables. Why is there a need for broadcast variables when working with Apache Spark? Explain with code snippet.</p> <p>b) How would you compute the total count of unique words in Spark? Write a program in Pyspark?</p>	1) The total number of movies	(2)	2) The maximum rating of movies	(2)	3) The number of movies that have maximum rating	(3)	4) The movies with ratings 1 and 2	(3)	5) The list of years and number of movies released each year	(5)	6) The number of movies that have a runtime of two hours	(5)	CO3
1) The total number of movies	(2)													
2) The maximum rating of movies	(2)													
3) The number of movies that have maximum rating	(3)													
4) The movies with ratings 1 and 2	(3)													
5) The list of years and number of movies released each year	(5)													
6) The number of movies that have a runtime of two hours	(5)													