



## END SEMESTER ASSESSMENT(ESA) B. TECH 5<sup>th</sup> Semester

### UE14CS314 - BIG DATA

Time: 2 Hrs

Answer All Questions

Max Marks: 60

Instructions:

- Provide short answers and to the point (max 3-4 lines) for all questions. Long answers when not asked for will not be looked upon favorably
- For problems, show the working. Providing just the answer is not acceptable.
- Answer all questions in the order they are asked. Leave space and come back if you need time to think
- One memory recall sheet (A4) is permitted in your own handwriting. Please attach that to your answer booklet

1		Provide short answers to the following questions. In case of True/False, please provide justification	10x2 =20
	a)	Which YARN scheduler uses a FIFO ordering for scheduling jobs?	
	b)	What security measures does Twitter take to deal with Photographs?	
	c)	What is the role of the SerDe component in HIVE?	
	d)	Why is a columnar database more suited for Analytics applications?	
	e)	What is a DRM in Mahout?	
	f)	If a Bloom filter determines that a key is in a set (e.g, a userid is in a set of known non-spam userids) then it is not necessary to look up the key from a backup store. True/False	
	g)	What are the function of a <i>supervisor</i> in Storm?	
	h)	What is the <i>batch size</i> in Streaming Spark?	
	i)	Which of the following operations is not supported by PNUTS? Why? Get, Set, Scan, Join	
	j)	How is load balancing achieved in PNUTS?	

2	a)	Suppose a stream of temperatures generated by weather monitoring stations is available. The records are in the form <time, location, temperature> which gives the time, location and value of temperature. For example, <2016_5_12.16:55:00, Bangalore-1, 25> means that at 16:55:00 on 12 May 2016, the temperature at the Bangalore-1 station was 25 degrees.  Draw the Storm topology needed for a Storm program that continuously outputs the value of the maximum temperature for each location for the day whenever it changes. Assume that the computation is too large to be handled by a single bolt. Show (i) how data is transferred between the various elements in the topology (ii) any data structures needed (iii) computation carried out in each bolt and output.	2+1+2
	a)	Explain the working of a Bloom filter with multiple hashes. Explain (i) data structures needed (ii) how the filter is initialized (iii) how it is used	2+1+2
	b)	Suppose the program of 2(a) is to be done using Streaming Spark and MapReduce. The	2+1+2

		<p>volume of the stream is too large to be handled by a single mapper or reducer.</p> <p>Draw the Streaming Spark workflow diagram needed for a program that continuously outputs the value of the current maximum temperature for each location for the day whenever it changes. Show (i) The incoming stream and conversion to an RDD stream (ii) What the mapper does (iii) any data structures needed (iv) computation carried out by the reducer(s).</p>	+2								
	c)	What is lineage in Spark? How is it used to handle failure of a task?	1+2								
3	a)	<p>Given that you have a HIVE installation with 2 tables in it -- <b>customers</b> and <b>purchases</b>. The customers table has the following structure &lt;cust_name, cust_id, dept&gt; and the purchases table has the structure &lt;purchase_id, cust_id, store, amount&gt; and you run the following SQL query on it:</p> <pre>select dept, count(1) from customers group by dept;</pre> <p>(i)What does this query do? (ii)How many map-reduce jobs will this require?</p>	2+3								
	b)	<p>Physically, HBase is composed of three types of servers in a master slave type of architecture.</p> <p>(i) What are the three types of servers?</p> <p>(ii) Why is Hbase a master-slave type of architecture?</p>	3+2								
	c)	<p>Cluster the following data into two clusters using a k-means map-reduce implementation. Assume that the first two points form the initial cluster centres. (i)Perform two iterations. For each iteration, show what the map and reduce stages will output. (ii) What will you use as a distance measure</p> <p>Data: 90, 80, 43, 26, 18, 75</p> <p>Express the solution in the following tabular format</p> <table border="1"> <thead> <tr> <th>Iteration 1 Map output</th><th>Iteration 1 Reduce Output</th><th>Iteration 2 Map Output</th><th>Iteration 2 Map Output</th></tr> </thead> <tbody> <tr> <td></td><td></td><td></td><td></td></tr> </tbody> </table>	Iteration 1 Map output	Iteration 1 Reduce Output	Iteration 2 Map Output	Iteration 2 Map Output					5
Iteration 1 Map output	Iteration 1 Reduce Output	Iteration 2 Map Output	Iteration 2 Map Output								
	d)	<p>"In YARN, Application Masters codify their request for resources using a ResourceRequest that contains – no of containers, resources per container, locality preferences and priority of requests within the application". (i)What is meant by resources per container? Why is this required?(ii)What are locality preferences?</p>	3+2								