

Student Answer Script View



MIT MPL-BTech-M Sc - MCA - 1st-3rd-5th and 7th Semester - Mid Term Examination - Sep 2024

Answer Sheet

Student Name: HARSH GONDAL . .

Roll Number: 220962334

Course: Computer Science and Engineering - Artificial Intelligence and Machine Learning

Year/Sem: Semester 5

Subject Name: BIG DATA ANALYTICS

Exam Date: 25-Sep-2024

Q.No : 1)

What does RDD stands for in Apache Spark?

**Resilient Distributed
Dataset**Reaction Distributed
Dataset

Reliable Distributed Dataset

Redundant Data Distribu

Q.No : 2)

Which of the following operations will cause an RDD to be recomputed in Spark?

count()

cache()

map()

filter()

Q.No : 3)

Which component of Hadoop is responsible for resource management?

☐ HBASE☒ YARN☐ Pig☐ Sqoop

Q.No : 4)

Which of the following is a NoSQL database in the Hadoop ecosystem?

HDFS

HIVE

HBASE

Sqoop

Q.No : 5)

Which of the following criteria is used by default to split nodes in PySpark's DecisionTreeClassifier?

Information Gain

Variance

Gini Index

Entropy

Q.No : 6)

Which method is used to make predictions using a trained decision tree model in PySpark?

transform()

fit()

predict()

forecast()

Q.No : 7)

Which of the following techniques is typically used to evaluate the performance of an ALS model in PySpark?

AUC

MSE

Precision

Recall

Q.No : 8)

For a data to be considered as Big Data following are necessary

Volume, Velocity, Value

Volume, Velocity, Variety

Volume, Velocity, Veracity

Volume, Veracity, Variety

Q.No : 9)

The method `wholeTextFile()` creates _____ RDD

☐ Double☒ Pair☐ Sequence☐ Union

Q.No : 10)

The function range() creates _____ RDD

Python

Sequence

Union

Parallel

Q.No : 11)

How do you schedule a YARN application? How YARN schedules and executes the application. Draw the YARN architecture.

Page:1

- The user either runs standalone mode or pseudo mode of Hadoop, which launches a single type for standalone mode, and JVM for every node in pseudo mode.
- This sends the application to the Resource Manager which is the master daemon of the yarn architecture. The Resource Manager directly responsible to client's request.
- The Resource Manager delegates a container to first or one Node Manager, this deleg process is called Application master. The Application master who figures out how much resources are required for the tasks and negotiate with Resource Manager to allocate resources to Node Managers, which run itself.
- The resources are then allocated by all other nodes, they starting processing.

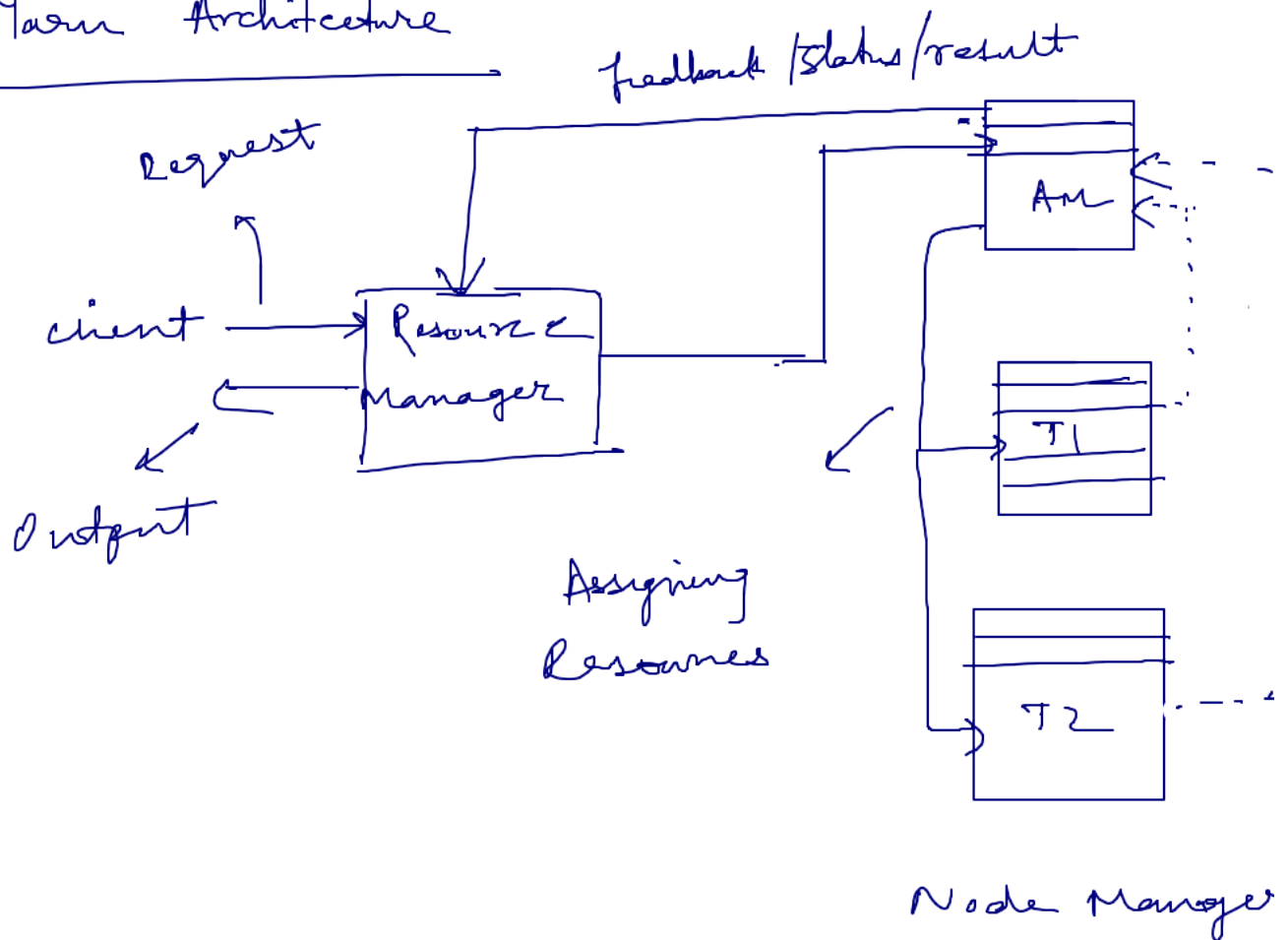
the data. The Application master keeps a ch

Page:2

of how much data is processed and feedbacks given time to time to Application master

- The Application Master on receiving feedback propagates it to Resource Manager which
- then to the client.

Yarn Architecture



Q.No : 12)

Write the code example for closure. How closure works?.

Page:1

Example of closure

```
def call_func (>):  
    def ret_message (concept):  
        return "This is a" + concept  
    return concept
```

closure is a python object in which scope is enclosed at the time the object is inst

So, when I call call_func (closure) , i will be printed this is a closure.

Q.No : 13)

Summarize different components of Hadoop Ecosystem.

Page:1

Different components of Hadoop Ecosystem are -

1) Mahout - A python library used for training machine learning algorithms

2) Pig - A high level language, which is used to analyze large dataset in its own language

3) Hive - It is a data warehousing system which facilitates storing data summarization ad hoc queries and analyzing large dataset. It has its own database structure language called HQL (Hive Query Language)

4) Hbase - It is a distributed, column oriented database, where data is stored in HDFS. HDFS supports both map reduce and queries.

5) Sqoop - It is a tool used for transferring bulk data into HDFS.

Page:2

6) Zoo keeper - It is a centralised service information configuration, naming and group various services

7) Ambari - It is a web based tool for provision, monitor of various hadoop clusters.

Q.No : 14)

Write a PySpark program to count the number of sentences in a text file. The sentences can be descriptive, interrogative, or exclamatory.

Page:1

```

from pyspark import SparkContext
sc = SparkContext()
sc.textFile(file_path).
  flatMap(lambda x: x.split('.')).
  flatMap(lambda x: x.split('?')).
  flatMap(lambda x: x.split('!')).
  filter(lambda x: x != '.' || x != '?' || x != '!')
sc.collect()

```


Q.No : 15)

How do you ensure that all the Hadoop environment variables are available for Hadoop installation? How do you create the index for HDFS? How do you check all Hadoop services are running?.

to ensure hadoop environment variables are available for hadoop installation, you can use the pseudo mode. In the pseudo mode JVM is created for all nodes ensures env variables are running.

The index of HDFS is created data is first loaded into the namenode.

The slave daemons keep updating their respective progress to the master node with checks they are running or not, where

the master nodes are checked, on the output received by the client and FsImage keep on updating.

Q.No : 16)

Write the code snippets for training the als model implicitly and explicitly.

Page:1

code snippets for ALS model are —

Explicitly—

```
als = ALS {  
    user : 'user_id'  
    item : 'item_id'  
    ratings : 'ratings_col' }  
als.fit(C)
```

Implicitly

```
als = ALS {  
    user : 'user_id'  
    item : 'item_id'  
}  
als.fit(C)
```


Q.No : 17)

Differentiate between top(5) and take(5) methods.

Page:1

- The top(5) method is a rdd method that collects the 5 elements from an RDD in decreasing order and prints it whereas take(5) collects the first 5 elements in an RDD and prints it.

Q.No : 18)

Summarize cogenerated RDD.

Page:1

Cogenerated RDD is a RDD formed by child and Parent RDD. Since RDD's are resilient that is there lineage is which refers to all the transformations in directed Acyclic graph. The Cogenerated is a partitioned according to its lineage.

Q.No : 19)

What is hyper parameter tuning? How do you achieve it in Decision trees? .

Page:1

Hyper parameter tuning refers to making grid and tune hyperparameter to machine learning models. Hyperparameters consists of functions sets constraints on training and evaluating model. It is done by `ParamGridBuilder()` which in conjunction with `CrossValidation()` or `TrainValidation` plots the grid.

In Decision Trees, the hyperparameters are

- max depth
- max bins
- Information Gain
- Impurity measure

So to tune these parameters, following code snippets-

from pyspark.ml.hyperparameter.tuning import
ParamGridBuilder

```
c = ParamGridBuilder()
```

```
c.addGrid('max-depth' : [4, 5]),
```

C. addGrid ('max-bins' : 20, 20, 1,

Page:2

c. addWord('impurity' : [a, b]),
c. addWord('regain' : [c, d]).

in such a way, the hyperparameters are tuned to particular which play a significant role, in training and evaluating model.

