Solution: Week 8 Assignment

Question 1: Research more on uber mode and create a detail.

Answer:

Uber Mode, the application is executed within the same container where the application master is running. The application master is responsible for coordinating the execution of tasks within a YARN application.

Uber Mode is primarily used for small jobs that don't require the resources of a full cluster. Instead of distributing the tasks across multiple containers in the cluster, Uber Mode allows the entire job to be executed within a single container, minimising the overhead associated with managing multiple containers and the associated inter-container communication.

Question 2: Explore a multimode cluster by navigating to the resource manager in the labs.

Answer:

The Resource Manager UI allows users to monitor and manage the resources and applications running on a YARN cluster. It provides a graphical interface to view various aspects of the cluster, including cluster utilisation, running applications, and resource allocation.

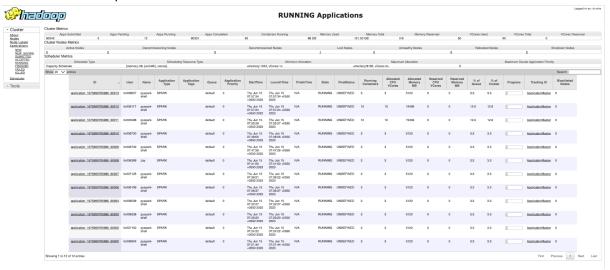
Here are some key features and components of the Resource Manager UI:

- Cluster Overview: The UI provides an overview of the cluster, including information such as the number of active nodes, total and available memory, CPU utilization, and the number of running and completed applications.
- 2. **Application List**: The UI displays a list of applications running on the cluster. It includes details like the application ID, user, state (e.g., running, completed, failed), start and finish time, and resource usage.
- 3. **Application Details**: Clicking on an application in the application list opens a detailed view. This view provides more information about the selected application, such as the list of containers allocated to the application, their resource usage, log links, and diagnostics information.
- 4. **Cluster Metrics**: The UI presents graphical representations of cluster metrics over time. These metrics may include CPU utilisation, memory usage, and other resource utilisation statistics.
- 5. **Queue Information**: If the cluster has multiple queues configured for resource allocation, the UI displays information about the queues. This includes details like the queue name, capacity, maximum capacity, current usage, and pending applications.
- 6. **Node Information**: The UI provides a summary of the nodes in the cluster. It shows details about each node, such as node ID, state, health status, available and used resources, and the number of containers running on the node.
- 7. **Logs and Diagnostics**: The Resource Manager UI allows you to access the logs and diagnostics information for individual applications. This can be helpful for troubleshooting issues or investigating the performance of specific applications.

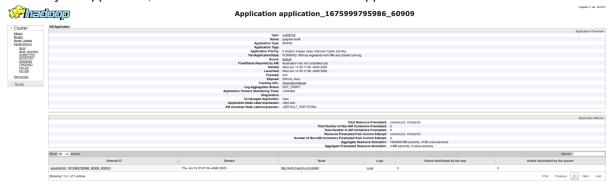
The Resource Manager UI is a powerful tool for monitoring and managing resources in a YARN cluster. It enables administrators and users to track the cluster's health, monitor application status, analyse resource usage, and diagnose problems. By leveraging the information provided by the UI, users can optimise resource allocation, identify bottlenecks, and ensure efficient cluster utilisation.

In the Resource Manager UI, you can see an overview of the applications, including the number of submitted and running applications. It also provides information about the memory and core usage out of the total available resources. Additionally, you can check the number of active nodes running in the cluster

As a user with the ID "itv006732," I am currently running a Spark session. I will explore the same information in the Spark UI.



Once you click on the application ID of the application you are running, you will be redirected to a page where you can access more detailed information about your application. This page provides information about the status of your application, the tracking URL for your job, the start and finish times of your application, and the current state of the YARN application.



Question4 : How to deal with nulls in Apache Spark? Please provide a detailed explanation with examples where necessary

Answer:

Dealing with null values in Apache Spark is an essential task to ensure accurate and reliable data processing. Apache Spark provides several methods and functions to handle null values effectively. In this explanation, I'll provide a detailed overview of common techniques and examples for dealing with nulls in Spark.

Dropping Null Values:

The simplest approach is to remove rows containing null values from your dataset. Spark provides the na object that allows you to perform operations on missing data. The drop() function can be used to drop rows with null values.

Example:

from pyspark.sql import SparkSession from pyspark.sql.functions import col spark = SparkSession.builder.getOrCreate() df = spark.read.option("header", "true").csv("data.csv") cleanedDF = df.na.drop() cleanedDF.show()

Filling Null Values:

Another approach is to fill the null values with specific default or derived values. Spark provides the fill() function to replace nulls with a specified value.

Example:

from pyspark.sql import SparkSession from pyspark.sql.functions import col spark = SparkSession.builder.getOrCreate() df = spark.read.option("header", "true").csv("data.csv") filledDF = df.na.fill("N/A") filledDF.show()

Additionally, you can fill nulls with different values based on the data type of the column. For numeric columns, you can use the fill() function with a numeric value, and for boolean columns, you can use fill() with a Boolean value.

Replacing Nulls with Conditional Values:

You may want to replace nulls with specific values based on certain conditions. Spark provides the when() and otherwise() functions, which allow you to define conditional logic while replacing nulls.

Example:

Filtering Null Values:

Sometimes, you may want to filter out rows containing null values for specific columns. You can use the filter() function along with isNull() or isNotNull() to achieve this.

Example:

```
from pyspark.sql import SparkSession
from pyspark.sql.functions import col
spark = SparkSession.builder.getOrCreate()
df = spark.read.option("header", "true").csv("data.csv")
filteredDF = df.filter(col("column1").isNotNull())
filteredDF.show()
```

Handling Nulls in Aggregations:

When performing aggregations on columns that contain nulls, Spark provides specialised functions such as coalesce(), first(), last(), and avg() that handle nulls appropriately.

Example:

These are some common techniques for dealing with null values in Apache Spark. The choice of method depends on the specific requirements of your data analysis or processing tasks. Remember to import the necessary functions from org.apache.spark.sql.functions to use the functions mentioned above.

Question 3: Select a relevant dataset of your choice and write some queries to demonstrate the following: running total, grouping aggregates, and various window functions like rank, dense_rank, row_number, lead, and lag. Also, try creating a pivot view.

Solution:

Link for code → https://g01.itversity.com/hub/user-redirect/lab/tree/g4_assignment.ipynb

Snapshot:

I am using hotel data to perform simple aggregations, such as counting. Here I am enforcing schema so we don't face any issue later with data type.

```
In [3]:

from pyspark.agl import SparkSession
import options
username = optipass.getuser()
spark = SparkSession.\
boilder.
boilder.
confid()*spark.agl.warehouse.dir', f'/user/itv008732/warehouse').\
eanbleNivoSupport().\
```

```
In [4]: hotel_df.show()
                                     customer_name|check_in_date|check_out_date|room_type|price|
                                      Jane Smith
                                                             2023-05-02
                                                                                   2023-05-06
                                                                                                      Deluxe 600.0
                                                                                   2023-05-08 Standard 450.0
2023-05-07 Executive 750.0
2023-05-09 Deluxe 550.0
                                                             2023-05-03
                                      Sarah Wilson
                                                             2023-05-04
                                                                                   2023-05-09 Deluxe 550.0
2023-05-10 Standard 400.0
                                    Emily Brown
                                                             2023-05-07
                              | Samantha Thompson
| William Lee
                                                             2023-05-08
                                                                                   2023-05-12 | Deluxe | 600.0
2023-05-13 | Standard | 450.0
                                 Amanda Harris
David Rodriguez
                                                             2023-05-11
2023-05-12
                                                                                   2023-05-16 | Executive | 750.0
2023-05-15 | Deluxe | 550.0
                                                                                   2023-05-16 | Executive | 550.0
2023-05-15 | Deluxe | 550.0
2023-05-18 | Standard | 400.0
                           10
                           11
                                      Linda Wilson
                                                             2023-05-14
                                                                                   2023-05-20 Deluxe | 600.0
2023-05-21 Standard | 450.0
2023-05-21 Standard | 450.0
2023-05-23 Executive | 750.0
2023-05-24 Deluxe | 550.0
                           12
                                   Robert Johnson
                                                             2023-05-15
                           13
                                 Sophia Anderson
                                                             2023-05-16
                                      James Smith
Olivia Brown
                                                             2023-05-17
2023-05-19
                           14
                           15
                                                                                   2023-05-25| Standard 400.0
                                   Michael Davis
Emily Thompson
William Lee
                          16
17
                                                             2023-05-20
                                                                                   2023-05-27
                                                                                                       Deluxe 600.0
                                                             2023-05-21
                                                                                   2023-05-28 Standard 450.0
2023-05-30 Executive 750.0
                           18
                                                             2023-05-23
                                         Ava Harris
                                                             2023-05-24
                           20
                                Daniel Rodriguez
                                                             2023-05-25
                                                                                   2023-05-29
                                                                                                       Deluxe|550.0
                                                                                   2023-06-01| Standard 400.0
                          21
                                    Sophia Wilson
                                                             2023-05-27
            only showing top 20 rows
In [5]: #let run the below command to use all sql functions in pyspark3
            from pyspark.sql.functions import *
```

Simple aggregations:

```
In [8]: # Create new session use hospital data set for the to perform more aggregation operation
  In [9]: from pyspark.sql import SparkSession
                  import getpass
                 import getpass
username = getpass.getuser()
spark = SparkSession. \
builder. \
config('spark.ui.port','0'). \
config('spark.sql.warehouse.dir', f'/user/itv006732/warehouse').\
enableHiveSupport(). \
matter('arr').
                 master('yarn'). \
getOrCreate()
In [10]: spark
Out[10]: SparkSession - hive
                  SparkContext
                  Spark UI
                  Version
                  v3.0.1
                  Master
                  varn
                  AppName
                  pyspark-shell
                  Grouping aggregations
In [11]:
                 hospital_df = spark.read \
    .format("csv") \
    .option("header", "true") \
    .option("inferschema", "true") \
    .load("/public/trendytech/datasets/hospital.csv")
                 hospital_df.printSchema()
                    root
|-- patient_id: integer (nullable = true)
|-- admission_date: string (nullable = true)
|-- discharge_date: string (nullable = true)
|-- diagnosis: string (nullable = true)
|-- doctor_id: integer (nullable = true)
|-- total_cost: double (nullable = true)
In [12]: hospital_df.show()
                   | \verb|patient_id|| admission_date|| discharge_date| \\ \qquad diagnosis|| doctor_id|| total_cost||
                                                                             2022-01-10 | Pneumonia |
2022-02-09 | Appendicitis |
2022-03-18 | Fractured Arm |
                                                 01-01-2022
                                                                                                                                         101
                                                                                                                                                        5000.0
                                                                                                                                         102
103
104
                                                 02-05-2022
                                                                                                                                                        7000.0
3500.0
                                                 03-12-2022
                                                                             2022-04-08 | Heart Attack |
2022-05-07 | Influenza |
2022-06-15 | Appendicitis |
2022-07-25 | Pneumonia |
                                                 04-02-2022
                                                                                                                                                      15000.0
                                                                                                                                         105
106
107
                                                 05-05-2022
                                                                                                                                                        2500.0
                                                06-10-2022
07-20-2022
                                                                                                                                                        5500.0
                                                                            2022-07-25 Pneumonia
2022-09-01 Heart Attack
2022-09-22 Fractured Leg
2022-10-10 Appendicitis
2022-12-18 Pneumonia
2022-12-18 Pneumonia
2023-01-09 Heart Attack
                                   8 |
9 |
10 |
                                                                                                                                         108
109
110
                                                08-25-2022
                                                                                                                                                      20000.0
                                                09-15-2022
10-05-2022
                                                                                                                                                       6000.0
7500.0
                                   11
                                                 11-02-2022
                                                                                                                                         111
                                                                                                                                                        2800.0
                                                12-10-2022
                                                                                                                                         112
                                                                                                                                                     18000.0
                                    13
                                                                          2023-01-09 Heart Attack
2023-02-18 Appendicitis
2023-03-28 Fractured Arm
2023-04-11 Influenza
2023-05-11 Heart Attack
                                                                                                                                                       7200.0
3800.0
2700.0
                                   14
                                                02-14-2023
                                                                                                                                         114
                                                 03-20-2023
                                                 04-05-2023
                                   16 |
17 |
                                                                                                                                         116
                                                05-08-2023
                                                                                                                                         117 İ
                                                                                                                                                     16000.0
```

```
In [13]: from pyspark.sql.functions import *
In [14]: #Let's convert the string to a datetime column for performing operations.
result_df1 = hospital_df.withColumn("admission_date", to_date("admission_date", "dd-mm-yyyy"))
result_df2 = result_df1.withColumn("discharge_date", to_date("discharge_date", "yyyy-mm-dd"))
In [15]: result_df2.show()
            |patient id|admission date|discharge date|
                                                                        diagnosis | doctor id | total cost |
                                 2022-01-01|
2022-01-02|
                                                     2022-01-10
                                                     2022-01-09 Appendicitis
                                                                                              102
                                                                                                         7000.0
                                 2022-01-03
                                                     2022-01-18 Fractured Arm
                                                                                              103
                                                                                                         3500.0
                                 2022-01-04
                                                    2022-01-10
2022-01-08
2022-01-07
                                                                                                       15000.0
                                                                    Heart Attack
                                                                        Influenza
                                                                                              105
                                                     2022-01-15 Appendicitis
                                 2022-01-06
                                                                                              106
                                                                                                        8000.0
                                                                                              107
108
109
                                                                    Pneumonia
Heart Attack
                                 2022-01-07
                                                     2022-01-25
                                                                                                         5500.0
                                                     2022-01-01
                                                                                                       20000.0
                                 2022-01-08
                                 2022-01-09
                                                     2022-01-22 Fractured Leg
                                                                                                         6000.0
                                                    2022-01-10 Appendicitis
2022-01-05 Influenza
2022-01-18 Pneumonia
                                                                                              110
111
112
                        10
                                 2022-01-10
                                                                                                         7500.0
                                 2022-01-11
                                                                                                         6000.0
                        12
                                                                                              113
114
115
                        13
                                 2023-01-01
                                                     2023-01-09
                                                                    Heart Attack
                                                                                                       18000.0
                        14
                                 2023-01-02
                                                     2023-01-18 | Appendicitis
2023-01-28 | Fractured Arm
                                                                                                         7200.0
                                                    2023-01-11 Influenza
2023-01-11 Heart Attack
2023-01-20 Pneumonia
                        16
                                 2023-01-04
                                                                                              116
                                                                                                         2700.0
                        17
                                 2023-01-05
                                                                                                        16000.0
                                 2023-01-06
                                                                                                         4800.0
                                                                        Pneumonia
                                                     2023-01-27 Fractured Leg
                        19
                                 2023-01-07
                                                                                              119
                                                                                                         6500.0
                        20
                                 2023-01-08
                                                    2023-01-16 Appendicitis
                                                                                                         7800.0
            only showing top 20 rows
In [16]: # Let's create a new column, 'year', which represents the year in which the patients are getting admitted."
            result_df3=result_df2.withColumn("year",year("admission_date"))
In [17]: result_df3.show()
            |patient id|admission date|discharge date|
                                                                        diagnosis | doctor id | total cost | year |
                                                    2022-01-10
                                 2022-01-01
                                                                        Pneumonia
                                                                                              101
                                                                                                         5000.0 | 2022
                                                                    Appendicitis
                                                                                                         7000.0 2022
                                 2022-01-02
                                                                                              102
                                 2022-01-03
                                                     2022-01-18 Fractured Arm
                                                                                              103
                                                                                                         3500.0 2022
                                                    2022-01-08
2022-01-07
                                                                    Heart Attack
                                                                                              104
                                 2022-01-04
                                                                                                       15000.0 2022
                                                                                                        2500.0|2022
8000.0|2022
                         6 |
7 |
                                                     2022-01-15 Appendicitis
                                 2022-01-06
                                                                                              106
                                 2022-01-07
                                                   2022-01-25
                                                                         Pneumonia
                                                                                              107 İ
                                                                                                        5500.0 2022
In [18]: # Let's find the number of patients diagnosed in each year, categorized by different categories.
summary_df=result_df3.groupby("year", "diagnosis").agg(count("patient_id").alias("total_patient")).sort(desc("year"))
summary_df.show()
            |year|
                       diagnosis|total_patient|
            2024
                       Influenza|
                       Influenza
            2023 Appendicitis
2023 Pneumonia
2023 Pneumonia
2023 Heart Attack
2023 Fractured Leg
2023 Fractured Arm
            2022 Appendicitis
            2022 Fractured Leg
            2022 Preumonia
2022 Preumonia
2022 Influenza
2022 Heart Attack
2022 Fractured Arm
In [19]: #For better visualization, let's create a pivot table based on this data and replace any null values with 0
           summary_df.groupBy("diagnosis").pivot("year").count().na.fill(0).show()
                 diagnosis|2022|2023|2024|
             Heart Attack
            Fractured Arm
                  Influenza
                 Pneumonia
```

```
In [20]: # Let's find the total number of patients for each category
          from pyspark.sql import Window
from pyspark.sql.functions import desc
          my_window=Window.partitionBy("diagnosis")
In [21]: summary_df2=result_df3.withColumn("total_patient_each_category",count("patient_id").over(my_window))
In [22]: summary_df2.show()
           | \texttt{patient\_id} | \texttt{admission\_date} | \texttt{discharge\_date} | \qquad \texttt{diagnosis} | \texttt{doctor\_id} | \texttt{total\_cost} | \texttt{year} | \texttt{total\_patient\_each\_category} | \\
                             2022-01-04|
                                             2022-01-08 | Heart Attack |
2022-01-01 | Heart Attack |
                                                                                 104
                                                                                         15000.0 | 2022 |
                             2022-01-08
                                                                                         20000.0 2022
                                                                                 108
                     13
                             2023-01-01
                                              2023-01-09
                                                           Heart Attack
                                                                                 113
                                                                                         18000.0 2023
                     17
                             2023-01-05
                                              2023-01-11
                                                                                         16000.0 2023
                                                           Heart Attack
                                                                                 117
                     22
                             2023-01-10
                                              2023-01-19 Heart Attack
                                                                                 122
                                                                                         21000.0 2023
                             2022-01-03
                                              2022-01-18 Fractured Arm
                                                                                 103
                                                                                          3500.0 2022
                     15
                             2023-01-03
                                              2023-01-28 Fractured Arm
                                                                                 115
                                                                                          3800.0 2023
                             2023-01-12
                                              2023-01-07 Fractured Arm
                             2022-01-09
                                              2022-01-22 Fractured Leg
                                                                                 109
                                                                                          6000.0 2022
                                                                                                                                      2
                     19
                                              2023-01-27 Fractured Leg
                                              2022-01-09 Appendicitis
2022-01-15 Appendicitis
                      2
                             2022-01-02
                                                                                 102
                                                                                           7000.0 2022
                             2022-01-06
                                                                                           8000.0 2022
                     10
                                              2022-01-10 Appendicitis
                             2022-01-10
                                                                                 110
                                                                                           7500.0 2022
                                                                                                                                      5
                     14
                             2023-01-02
                                              2023-01-18
                                                           Appendicitis
                                                                                           7200.0 2023
                     20
                             2023-01-08
                                              2023-01-16 Appendicitis
                                                                                 120
                                                                                           7800.0 2023
                                                                                                                                      5
                     5 |
                             2022-01-05
                                              2022-01-07
                                                               Influenza
                                                                                 105
                                                                                           2500.0 2022
                             2022-01-11
                                              2022-01-05
                                                                                 111
                                                                                          2800.0 2022
                                                               Influenza
                                                                                                                                      5
                     16
                             2023-01-04
                                              2023-01-11
                                                               Influenza
                                                                                           2700.0 2023
                             2023-01-09
                                              2023-01-09
                                                               Influenza
                                                                                           2900.0 2023
                                                                                 121
                     25
                             2024-01-01
                                              2024-01-15
                                                               Influenza
                                                                                          3200.0 2024
                                                                                                                                      5
          only showing top 20 rows
In [231:
```

```
In [23]:
    from pyspark.sql.functions import row_number
    from pyspark.sql.window import Window

my_window2 = Window.partitionBy("diagnosis").orderBy(summary_df2["total_patient_each_category"].desc())

summary_df3 = summary_df2.withColumn("row_num", row_number().over(my_window2))

top3_df = summary_df3.filter(summary_df3["row_num"] <= 1)

#let's take neccessery column

top3_df['diagnosis','total_patient_each_category'].show()</pre>
```

```
In [24]: #Let's find the top 3 patients who have the highest total cost for each diagnosis.
               my_window3=Window.partitionBy("diagnosis").orderBy(desc("total_cost"))
summary_df4=result_df3.withColumn("rank",rank().over(my_window3))
top3_df=summary_df4.filter(summary_df4["rank"]<=3)</pre>
In [25]: top3_df.show()
                | \verb|patient_id|| \verb|admission_date|| \verb|discharge_date|| \qquad | \verb|diagnosis|| \verb|doctor_id|| total_cost|| year|| rank||
                                                                  2023-01-19 | Heart Attack
2022-01-01 | Heart Attack
2023-01-09 | Heart Attack
2023-01-07 | Fractured Arm
2023-01-28 | Fractured Arm
                                                                                                                                 21000.0|2023|
                                          2023-01-10
                                          2022-01-08
                                                                                                                      108
                                                                                                                                 20000.0 2022
                                         2023-01-01
2023-01-12
2023-01-03
                                                                                                                                 18000.0|2023|
4100.0|2023|
3800.0|2023|
                                                                                                                      113
124
                              13
                              15
                                                                                                                      115
                                                                  2022-01-28 Fractured Arm
2023-01-27 Fractured Leg
2022-01-22 Fractured Leg
                                         2022-01-03
                                                                                                                      103
                                                                                                                                   3500.0 2022
                                         2023-01-07
2022-01-09
                                                                                                                      119
109
                                                                                                                                   6500.0|2023
6000.0|2022
                                                                  2022-01-15 | Appendicitis
2023-01-16 | Appendicitis
2022-01-10 | Appendicitis
2024-01-15 | Influenza
                                         2022-01-06
                                                                                                                      106
                                                                                                                                   8000.0 2022
                                         2023-01-08
2023-01-10
2022-01-10
2024-01-01
                              20
10
                                                                                                                      120
110
                                                                                                                                   7800.0 2023
7500.0 2022
                                                                                                                                   3200.0 2024
                              25
                                                                                                                      125
                              21
                                          2023-01-09
                                                                  2023-01-09
                                                                                           Influenza
                                                                                                                      121
                                                                                                                                   2900.0 2023
                                          2022-01-11
                                                                  2022-01-05
2022-01-18
                                                                                                                      111
112
                                                                                                                                   2800.0 2022
                              12
                                                                                                                                   6000.0 2022
                                                                                           Pneumonia
                                          2022-01-07
                                                                  2022-01-25
                                                                                           Pneumonia
                                                                                                                      107
                                                                                                                                   5500.0 2022
                              23
                                          2023-01-11
                                                                  2023-01-22
                                                                                           Pneumonia
                                                                                                                      123
                                                                                                                                   5200.0 2023
In [26]: # let's drop the rank column there is no use as of now for this.
               top3_df.drop("rank","year").show()
                |patient_id|admission_date|discharge_date| diagnosis|doctor_id|total_cost|
                                                                 2023-01-19 | Heart Attack
2022-01-01 | Heart Attack
2023-01-09 | Heart Attack
2023-01-07 | Fractured Arm
                                         2023-01-10|
                                                                                                                                 21000.0
                              22 |
                                                                                                                      122
                                         2022-01-08
2023-01-01
2023-01-12
                                                                                                                      108
113
                                                                                                                                 20000.0
                              24
                                                                                                                      124
                                                                                                                                   4100.0
                                                                  2023-01-28 Fractured Arm
2022-01-18 Fractured Arm
2023-01-27 Fractured Leg
2022-01-22 Fractured Leg
                                                                                                                      115
103
119
                              15
                                          2023-01-03
                                                                                                                                   3800.0
                                         2022-01-03
2023-01-07
                                                                                                                                   3500.0
                              19
                                                                                                                                   6500.0
                                         2022-01-09
                                                                                                                      109
                                                                                                                                   6000.0
```

8000.0 7800.0 7500.0

106 120

2022-01-06 2022-01-08 2023-01-10 2022-01-15| Appendicitis| 2023-01-16| Appendicitis| 2023-01-10| Appendicitis|