# QUESTION 1

**What is UBER mode ?**
When a particular Hadoop job is invoked from client machine, RM will create separate container for that application. Uber configuration, will allow running that job in the same process as the Application Master (AM).

**Uber jobs:**
Uber jobs are jobs that are executed within the ApplicationMaster. It do not need to communicate with RM to create the containers on different worker nodes. As the job is given is evaluated small so the AM runs tasks within its own process and avoided the overhead of launching and communicate with remote containers located in different worker nodes can be considered as advantage in certain cases.
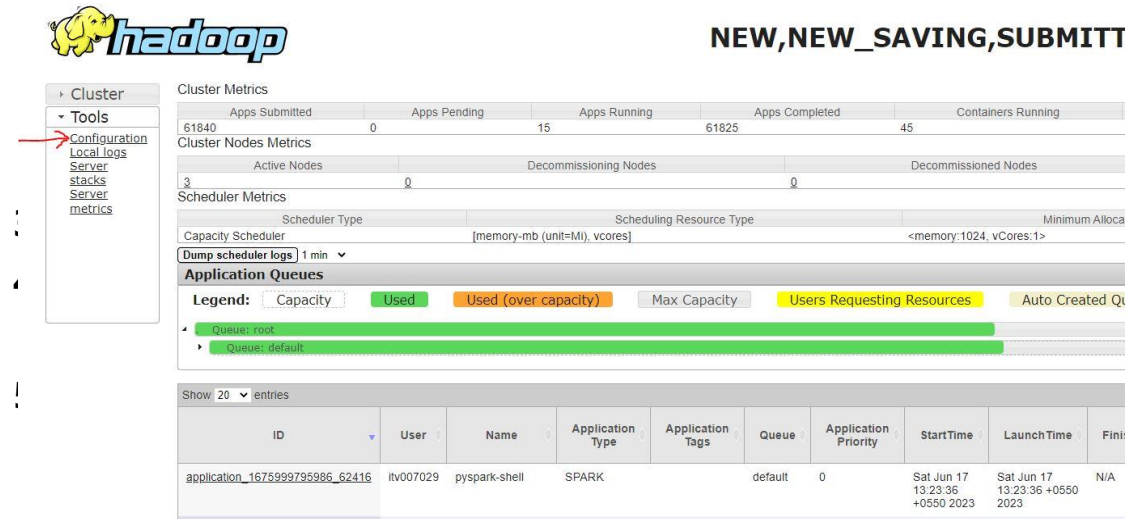
**Why?**
If you have a small dataset, Uber configuration will help you out, by reducing additional time required to complete the certain Hadoop job.

## Configuration options for Uber Jobs

| Property | Description |
|---|---|
| mapreduce.job.ubertask.enable | Whether to enable the small-jobs "ubertask" optimization, which runs "sufficiently small" jobs sequentially within a single JVM. "Small" is defined by the following maxmaps, maxreduces, and maxbytes settings. Users can override this value.<br><br>Default = false |
| mapreduce.job.ubertask.maxmaps | The threshold for the number of maps beyond which a job is considered too large for the ubertasking optimization. Users can override this value, but only downward.<br><br>Default = 9 |
| mapreduce.job.ubertask.maxreduces | The threshold for the number of reduces beyond which a job is considered too large for the uber tasking optimization. CURRENTLY THE CODE CANNOT SUPPORT MORE THAN ONE REDUCE and will ignore larger values (zero is a valid maximum value, however). Users can override this value, but only downward.<br><br>Default = 1 |
| mapreduce.job.ubertask.maxbytes | The threshold for the number of input bytes beyond which a job is considered too large for the uber tasking optimization. If no value is specified, dfs.block.size is used as the default. Be sure to specify a default value in mapred-site.xml if the underlying file system is not HDFS. Users can override this value, but only downward.<br><br>Default = HDFS Block Size |

# SOP for accessing UBER mode in Resource Manager –

1. Go to http://172.16.1.104:19088/ by enabling foxy proxy server.
2. Click on Configuration option under Tool tab on left side of screen.



3. Find in page for the keyword "Uber". Setting will appear change as mentioned.

# QUESTION 2

**ABOUT** - Gives the information about the cluster and it's metrics eg – memory used, total memory, v cores etc. Other information like Cluster Node Metrics and Scheduler Metrics (min and max vcore allocation info). We can also find certain details related to commissioning of Resource Manager.

**About the Cluster**

Cluster Metrics

| Apps Submitted | Apps Pending | Apps Running | Apps Completed | Containers Running | Memory Used | Memory Total | Memory Reserved | VCores Used | VCores Total | VCores Reserved |
|---|---|---|---|---|---|---|---|---|---|---|
| 61881 | 0 | 14 | 61867 | 42 | 70 GB | 151.00 GB | 0 B | 42 | 90 | 0 |

Cluster Nodes Metrics

| Active Nodes | Decommissioning Nodes | Decommissioned Nodes | Lost Nodes | Unhealthy Nodes | Rebooted Nodes | Shutdown Nodes |
|---|---|---|---|---|---|---|
| 3 | 0 | 0 | 1 | 0 | 0 | 0 |

Scheduler Metrics

| Scheduler Type | Scheduling Resource Type | Minimum Allocation | Maximum Allocation | Maximum Cluster Application Priority |
|---|---|---|---|---|
| Capacity Scheduler | [memory-mb (unit=Mi), vcores] | <memory:1024, vCores:1> | <memory:8192, vCores:4> | 0 |

Cluster overview

| | |
|---|---|
| Cluster ID: | 1675999795986 |
| ResourceManager state: | STARTED |
| ResourceManager HA state: | active |
| ResourceManager HA zookeeper connection state: | Could not find leader elector. Verify both HA and automatic failover are enabled. |
| ResourceManager RMStateStore: | org.apache.hadoop.yarn.server.resourcemanager.recovery.NullRMStateStore |
| ResourceManager started on: | Thu Feb 09 22:29:55 -0500 2023 |
| ResourceManager version: | 3.3.0 from aa96f1871bfd858f9bac59cf2a81ec470da649af by brahma source checksum e0a276649f889c15d0e8f08eccd0c10 on 2020-07-06T18:58Z |
| Hadoop version: | 3.3.0 from aa96f1871bfd858f9bac59cf2a81ec470da649af by brahma source checksum 5dc29b802d6ccd77b262ef9d04d19c4 on 2020-07-06T18:44Z |

**Nodes** – Shows information about 3 worker nodes present on this cluster.

**Nodes of the cluster**

Cluster Metrics

| Apps Submitted | Apps Pending | Apps Running | Apps Completed | Containers Running | Memory Used | Memory Total | Memory Reserved | VCores Used | VCores Total | VCores Reserved |
|---|---|---|---|---|---|---|---|---|---|---|
| 61888 | 0 | 15 | 61873 | 45 | 75 GB | 151.00 GB | 0 B | 45 | 90 | 0 |

Cluster Nodes Metrics

| Active Nodes | Decommissioning Nodes | Decommissioned Nodes | Lost Nodes | Unhealthy Nodes | Rebooted Nodes | Shutdown Nodes |
|---|---|---|---|---|---|---|
| 3 | 0 | 0 | 1 | 0 | 0 | 0 |

Scheduler Metrics

| Scheduler Type | Scheduling Resource Type | Minimum Allocation | Maximum Allocation | Maximum Cluster Application Priority |
|---|---|---|---|---|
| Capacity Scheduler | [memory-mb (unit=Mi), vcores] | <memory:1024, vCores:1> | <memory:8192, vCores:4> | 0 |

Show 20 entries                                                                 Search:

| Node Labels | Rack | Node State | Node Address | Node HTTP Address | Last health-update | Health-report | Containers | Allocation Tags | Mem Used | Mem Avail | VCores Used | VCores Avail | Version |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | /default-rack | RUNNING | w01.itversity.com:35127 | w01.itversity.com:8042 | Sat Jun 17 05:15:41 -0400 2023 | | 10 | | 15 GB | 35.33 GB | 10 | 20 | 3.3.0 |
| | /default-rack | RUNNING | w03.itversity.com:41791 | w03.itversity.com:8042 | Sat Jun 17 05:14:55 -0400 2023 | | 20 | | 34 GB | 16.33 GB | 20 | 10 | 3.3.0 |
| | /default-rack | RUNNING | w02.itversity.com:46669 | w02.itversity.com:8042 | Sat Jun 17 05:16:29 -0400 2023 | | 15 | | 26 GB | 24.33 GB | 15 | 15 | 3.3.0 |

Showing 1 to 3 of 3 entries                    First    Previous    1    Next    Last

By clicking on link provided under Node HTTP Address of particular worker node we can get info about that particular node and as well **node manager** information.

NodeManager information

| | |
|---|---|
| Total Vmem allocated for Containers | 251.66 GB |
| Vmem enforcement enabled | true |
| Total Pmem allocated for Container | 50.33 GB |
| Pmem enforcement enabled | true |
| Total VCores allocated for Containers | 30 |
| Resource types | memory-mb (unit=Mi), vcores |
| NodeHealthy Status | true |
| LastNodeHealthTime | Sat Jun 17 05:35:41 EDT 2023 |
| NodeHealthReport | |
| NodeManager started on | Wed Feb 15 08:19:38 EST 2023 |
| NodeManager Version | 3.3.0 from aa96f1871bfd858f9bac59cf2a81ec470da649af by brahma source checksum e0a276649f889c15d0e8f08eccd0c10 on 2020-07-06T18:58Z |
| Hadoop Version | 3.3.0 from aa96f1871bfd858f9bac59cf2a81ec470da649af by brahma source checksum 5dc29b802d6ccd77b262ef9d04d19c4 on 2020-07-06T18:44Z |

**Applications** – Shows all the links and key information about all the application present on this cluster, where different tab under applications can show particular job which are running, finished or killed etc.



Link under the ID of particular running application will lead to that particular application information in detailed view.



By clicking on link under Tracking URL named as Application Master will lead us to **Spark UI** of that particular job and we can navigate through various information about the job in that particular application.

**Link under Description** section will lead us to DAG Visualization and stages information. **Stages tab** show completed stages. Storage tab show the cached results. **Environment tab** gives us information about spark properties, system properties etc. **Executors tab** show the information about information about driver and executors involved for this application in order to complete a job. When clicked on **SQL tab** shows the completed queries and when clicked in description details of the query with physical, optimized, analyzed and parsed logical plan.

**TOOLS** – under tools section we can open configuration tab in order to edit predefined configurations and property related to resource manager e.g. – changing uber mode configuration.

# QUESTION 3

```
[1]: from pyspark.sql import SparkSession
     import getpass
     username = getpass.getuser()
     spark = SparkSession. \
     builder. \
     appName("Ultron"). \
     config('spark.ui.port','0'). \
     config("spark.sql.warehouse.dir", f"/user/itv005937/warehouse"). \
     enableHiveSupport(). \
     master('yarn'). \
     getOrCreate()
```

```
[2]: groceries_schema = "order_id integer, location string, item string, order_date string, quantity integer"
```

```
[3]: from pyspark.sql import *
     from pyspark.sql.functions import *
```

```
[4]: groceries_new_df = spark.read \
     .format("csv") \
     .schema(groceries_schema) \
     .option("header","true") \
     .load("/user/itv005937/data/groceries_new.csv")
```

```
[5]: groceries_df = groceries_new_df.withColumn("order_Date", to_date("order_date","dd/mm/yyyy"))
```

```
[6]: groceries_df.createOrReplaceTempView("groceries_sql")
```

## Aggregate function

```
[7]: groceries_df.select(count("*").alias("No_of_Records"),
                         countDistinct("location").alias("No_of_location"),
                         sum("quantity").alias("Total_Quantity")).show()
```

```
+-------------+--------------+--------------+
|No_of_Records|No_of_location|Total_Quantity|
+-------------+--------------+--------------+
|           21|             7|           273|
+-------------+--------------+--------------+
```

```
[8]: groceries_df.selectExpr("count(*) as No_of_Records",
                            "count(distinct(location)) as No_of_location",
                            "sum(quantity) as Total_Quantity").show()
```

```
+-------------+--------------+--------------+
|No_of_Records|No_of_location|Total_Quantity|
+-------------+--------------+--------------+
|           21|             7|           273|
+-------------+--------------+--------------+
```

```
                            count(distinct(location)) as No_of_location ,
                            "sum(quantity) as Total_Quantity").show()
```

```
+-------------+--------------+--------------+
|No_of_Records|No_of_location|Total_Quantity|
+-------------+--------------+--------------+
|           21|             7|           273|
+-------------+--------------+--------------+
```

[9]:
```
spark.sql("""select count(*) as No_of_Records,
count(distinct(location)) as No_of_location,
sum(quantity) as Total_Quantity from groceries_sql""").show()
```

```
+-------------+--------------+--------------+
|No_of_Records|No_of_location|Total_Quantity|
+-------------+--------------+--------------+
|           21|             7|           273|
+-------------+--------------+--------------+
```

## Grouping Aggregation

[10]:
```
GA1 = groceries_df.groupBy("location","item") \
.agg(sum("quantity").alias("Total_qunatity_per_group"), sum(expr("(quantity/273)*100")) \
.alias("percentage_of_total_quantity")).sort("location")
```

[11]: `GA1.show()`

```
+---------+--------+------------------------+----------------------------+
| location|    item|Total_qunatity_per_group|percentage_of_total_quantity|
+---------+--------+------------------------+----------------------------+
| Bellevue| Flowers|                      10|           3.6630036630036633|
| Bellevue|   Bread|                     125|           45.78754578754578 |
| Issaquah|    Meat|                      14|           5.128205128205129 |
| Issaquah|   Onion|                      12|           4.395604395604395 |
| Issaquah|  Tomato|                       6|           2.197802197802198 |
|     Kent|  Apples|                      20|           7.326007326007327 |
|  Redmond|  Cheese|                      15|           5.494505494505495 |
|  Redmond|    Meat|                      40|          14.652014652014653 |
|  Redmond|   Bread|                       5|           1.8315018315018317|
|   Renton|   Bread|                       5|           1.8315018315018317|
|Sammanish|   Bread|                       5|           1.8315018315018317|
|  Seattle| Bananas|                       7|           2.564102564102564 |
|  Seattle|Potatoes|                       9|           3.296703296703297 |
+---------+--------+------------------------+----------------------------+
```

[12]:
```
GA2 = groceries_df.groupBy("location","item") \
.agg(expr("sum(quantity) as Total_qunatity_per_group"), expr("sum((quantity/273)*100) as percentage_of_total_quantity")) \
.sort("location")
```

[13]: `GA2.show()`

```
+---------+--------+------------------------+----------------------------+
| location|    item|Total_qunatity_per_group|percentage_of_total_quantity|
+---------+--------+------------------------+----------------------------+
| Bellevue| Flowers|                      10|           3.6630036630036633|
| Bellevue|   Bread|                     125|           45.78754578754578 |
| Issaquah|    Meat|                      14|           5.128205128205129 |
| Issaquah|   Onion|                      12|           4.395604395604395 |
| Issaquah|  Tomato|                       6|           2.197802197802198 |
|     Kent|  Apples|                      20|           7.326007326007327 |
|  Redmond|  Cheese|                      15|           5.494505494505495 |
|  Redmond|   Bread|                       5|           1.8315018315018317|
|  Redmond|    Meat|                      40|          14.652014652014653 |
|   Renton|   Bread|                       5|           1.8315018315018317|
|Sammanish|   Bread|                       5|           1.8315018315018317|
|  Seattle| Bananas|                       7|           2.564102564102564 |
|  Seattle|Potatoes|                       9|           3.296703296703297 |
+---------+--------+------------------------+----------------------------+
```

[15]:
```
spark.sql("""select location, item, sum(quantity) as Total_qunatity_per_group,
sum((quantity/273)*100) as percentage_of_total_quantity
from groceries_sql group by location,item order by location""").show()
```

```
+---------+--------+------------------------+----------------------------+
| location|    item|Total_qunatity_per_group|percentage_of_total_quantity|
+---------+--------+------------------------+----------------------------+
| Bellevue| Flowers|                      10|           3.6630036630036633|
| Bellevue|   Bread|                     125|           45.78754578754578 |
| Issaquah|    Meat|                      14|           5.128205128205129 |
| Issaquah|   Onion|                      12|           4.395604395604395 |
| Issaquah|  Tomato|                       6|           2.197802197802198 |
|     Kent|  Apples|                      20|           7.326007326007327 |
|  Redmond|   Bread|                       5|           1.8315018315018317|
|  Redmond|  Cheese|                      15|           5.494505494505495 |
|  Redmond|    Meat|                      40|          14.65201465201465  |
|   Renton|   Bread|                       5|           1.8315018315018317|
|Sammanish|   Bread|                       5|           1.8315018315018317|
|  Seattle| Bananas|                       7|           2.564102564102564 |
|  Seattle|Potatoes|                       9|           3.296703296703297 |
+---------+--------+------------------------+----------------------------+
```

# Window Functions - Running Total, Rank, Dense Rank, Row Num, Lead, Lag

```
[9]: window_open_1 = Window.partitionBy("location","item") \
     .orderBy("order_date") \
     .rowsBetween(Window.unboundedPreceding, Window.currentRow)
```

```
[10]: RT_df = groceries_df.withColumn("Running_Total", sum("quantity").over(window_open_1))
```

```
[11]: RT_df.show()
```

```
+--------+---------+--------+----------+--------+-------------+
|order_id| location|    item|order_Date|quantity|Running_Total|
+--------+---------+--------+----------+--------+-------------+
|       2|     Kent|  Apples|2017-01-02|      20|           20|
|      13|Sammamish|   Bread|2017-01-07|       5|            5|
|      15| Issaquah|    Meat|2017-01-08|       3|            3|
|      16| Issaquah|    Meat|2017-01-09|       5|            8|
|      17| Issaquah|    Meat|2017-01-10|       6|           14|
|      11|   Renton|   Bread|2017-01-05|       5|            5|
|       1|  Seattle| Bananas|2017-01-01|       7|            7|
|       5|  Seattle|Potatoes|2017-01-04|       9|            9|
|       9|  Redmond|  Cheese|2017-01-05|      15|           15|
|       8| Issaquah|   Onion|2017-01-05|       4|            4|
|      10| Issaquah|   Onion|2017-01-06|       4|            8|
|      12| Issaquah|   Onion|2017-01-07|       4|           12|
|       7|  Redmond|   Bread|2017-01-05|       5|            5|
|       4|  Redmond|    Meat|2017-01-03|      40|           40|
|       3| Bellevue| Flowers|2017-01-02|      10|           10|
|      14| Issaquah|   Tomato|2017-01-07|      6|            6|
|       6| Bellevue|   Bread|2017-01-04|       5|            5|
|      18| Bellevue|   Bread|2017-01-11|       7|           12|
|      19| Bellevue|   Bread|2017-01-12|      54|           66|
|      20| Bellevue|   Bread|2017-01-13|      34|          100|
+--------+---------+--------+----------+--------+-------------+
only showing top 20 rows
```

```
[12]: window_open_2 = Window.partitionBy("location") \
     .orderBy("order_date") \
```

```
[13]: Rank_df = groceries_df.withColumn("Rank", rank().over(window_open_2))
```

```
[14]: Rank_df.show()
```

```
+--------+---------+--------+----------+--------+----+
|order_id| location|    item|order_Date|quantity|Rank|
+--------+---------+--------+----------+--------+----+
|       8| Issaquah|   Onion|2017-01-05|       4|   1|
|      10| Issaquah|   Onion|2017-01-06|       4|   2|
|      12| Issaquah|   Onion|2017-01-07|       4|   3|
|      14| Issaquah|  Tomato|2017-01-07|       6|   3|
|      15| Issaquah|    Meat|2017-01-08|       3|   5|
|      16| Issaquah|    Meat|2017-01-09|       5|   6|
|      17| Issaquah|    Meat|2017-01-10|       6|   7|
|      13|Sammamish|   Bread|2017-01-07|       5|   1|
|       4|  Redmond|    Meat|2017-01-03|      40|   1|
|       7|  Redmond|   Bread|2017-01-05|       5|   2|
|       9|  Redmond|  Cheese|2017-01-05|      15|   2|
|       1|  Seattle| Bananas|2017-01-01|       7|   1|
|       5|  Seattle|Potatoes|2017-01-04|       9|   2|
|       2|     Kent|  Apples|2017-01-02|      20|   1|
|       3| Bellevue| Flowers|2017-01-02|      10|   1|
|       6| Bellevue|   Bread|2017-01-04|       5|   2|
|      18| Bellevue|   Bread|2017-01-11|       7|   3|
|      19| Bellevue|   Bread|2017-01-12|      54|   4|
|      20| Bellevue|   Bread|2017-01-13|      34|   5|
|      21| Bellevue|   Bread|2017-01-14|      25|   6|
+--------+---------+--------+----------+--------+----+
only showing top 20 rows
```

```
[15]: Dense_Rank_df = groceries_df.withColumn("Dense_Rank", dense_rank().over(window_open_2))
```

```
[16]: Dense_Rank_df.show()
```

```
+--------+---------+-------+----------+--------+----------+
|order_id| location|   item|order_Date|quantity|Dense_Rank|
+--------+---------+-------+----------+--------+----------+
|       8| Issaquah|  Onion|2017-01-05|       4|         1|
|      10| Issaquah|  Onion|2017-01-06|       4|         2|
|      12| Issaquah|  Onion|2017-01-07|       4|         3|
|      14| Issaquah| Tomato|2017-01-07|       6|         3|
|      15| Issaquah|   Meat|2017-01-08|       3|         4|
|      16| Issaquah|   Meat|2017-01-09|       5|         5|
|      17| Issaquah|   Meat|2017-01-10|       6|         6|
|      13|Sammamish|  Bread|2017-01-07|       5|         1|
|       4|  Redmond|   Meat|2017-01-03|      40|         1|
|       7|  Redmond|  Bread|2017-01-05|       5|         2|
|       9|  Redmond| Cheese|2017-01-05|      15|         2|
|       1|  Seattle|Bananas|2017-01-01|       7|         1|
|       5|  Seattle|Potatoes|2017-01-04|      9|         2|
|       2|     Kent| Apples|2017-01-02|      20|         1|
|       3| Bellevue|Flowers|2017-01-02|      10|         1|
|       6| Bellevue|  Bread|2017-01-04|       5|         2|
|      18| Bellevue|  Bread|2017-01-11|       7|         3|
|      19| Bellevue|  Bread|2017-01-12|      54|         4|
|      20| Bellevue|  Bread|2017-01-13|      34|         5|
|      21| Bellevue|  Bread|2017-01-14|      25|         6|
+--------+---------+-------+----------+--------+----------+
only showing top 20 rows
```

```
[17]: Row_Num_df = groceries_df.withColumn("Row_Num", row_number().over(window_open_2))
```

```
[18]: Row_Num_df.show()
```

```
+--------+---------+-------+----------+--------+-------+
|order_id| location|   item|order_Date|quantity|Row_Num|
+--------+---------+-------+----------+--------+-------+
|       8| Issaquah|  Onion|2017-01-05|       4|      1|
|      10| Issaquah|  Onion|2017-01-06|       4|      2|
|      12| Issaquah|  Onion|2017-01-07|       4|      3|
|      14| Issaquah| Tomato|2017-01-07|       6|      4|
|      15| Issaquah|   Meat|2017-01-08|       3|      5|
|      16| Issaquah|   Meat|2017-01-09|       5|      6|
|      17| Issaquah|   Meat|2017-01-10|       6|      7|
|      13|Sammamish|  Bread|2017-01-07|       5|      1|
|       4|  Redmond|   Meat|2017-01-03|      40|      1|
|       7|  Redmond|  Bread|2017-01-05|       5|      2|
|       9|  Redmond| Cheese|2017-01-05|      15|      3|
|       1|  Seattle|Bananas|2017-01-01|       7|      1|
|       5|  Seattle|Potatoes|2017-01-04|      9|      2|
|       2|     Kent| Apples|2017-01-02|      20|      1|
|       3| Bellevue|Flowers|2017-01-02|      10|      1|
|       6| Bellevue|  Bread|2017-01-04|       5|      2|
|      18| Bellevue|  Bread|2017-01-11|       7|      3|
|      19| Bellevue|  Bread|2017-01-12|      54|      4|
|      20| Bellevue|  Bread|2017-01-13|      34|      5|
|      21| Bellevue|  Bread|2017-01-14|      25|      6|
+--------+---------+-------+----------+--------+-------+
only showing top 20 rows
```

```
[19]: Lead_df = groceries_df.withColumn("Upcoming_quantity", lead("quantity").over(window_open_2))
```

```
[20]: Lead_df.show()
```

```
+--------+---------+-------+----------+--------+-----------------+
|order_id| location|   item|order_Date|quantity|Upcoming_quantity|
+--------+---------+-------+----------+--------+-----------------+
|       8| Issaquah|  Onion|2017-01-05|       4|                4|
|      10| Issaquah|  Onion|2017-01-06|       4|                4|
|      12| Issaquah|  Onion|2017-01-07|       4|                6|
|      14| Issaquah| Tomato|2017-01-07|       6|                3|
|      15| Issaquah|   Meat|2017-01-08|       3|                5|
|      16| Issaquah|   Meat|2017-01-09|       5|                6|
|      17| Issaquah|   Meat|2017-01-10|       6|             null|
|      13|Sammamish|  Bread|2017-01-07|       5|             null|
|       4|  Redmond|   Meat|2017-01-03|      40|                5|
|       7|  Redmond|  Bread|2017-01-05|       5|               15|
|       9|  Redmond| Cheese|2017-01-05|      15|             null|
|       1|  Seattle|Bananas|2017-01-01|       7|                9|
|       5|  Seattle|Potatoes|2017-01-04|      9|             null|
|       2|     Kent| Apples|2017-01-02|      20|             null|
|       3| Bellevue|Flowers|2017-01-02|      10|                5|
|       6| Bellevue|  Bread|2017-01-04|       5|                7|
|      18| Bellevue|  Bread|2017-01-11|       7|               54|
|      19| Bellevue|  Bread|2017-01-12|      54|               34|
|      20| Bellevue|  Bread|2017-01-13|      34|               25|
|      21| Bellevue|  Bread|2017-01-14|      25|             null|
+--------+---------+-------+----------+--------+-----------------+
only showing top 20 rows
```

```
[21]: Lag_df = groceries_df.withColumn("Previous_quantity", lag("quantity").over(window_open_2))
```

```
[22]: Lag_df.show()
```

```
+--------+---------+-------+----------+--------+-----------------+
|order_id| location|   item|order_Date|quantity|Previous_quantity|
+--------+---------+-------+----------+--------+-----------------+
|       8| Issaquah|  Onion|2017-01-05|       4|             null|
|      10| Issaquah|  Onion|2017-01-06|       4|                4|
|      12| Issaquah|  Onion|2017-01-07|       4|                4|
|      14| Issaquah| Tomato|2017-01-07|       6|                4|
|      15| Issaquah|   Meat|2017-01-08|       3|                6|
|      16| Issaquah|   Meat|2017-01-09|       5|                3|
|      17| Issaquah|   Meat|2017-01-10|       6|                5|
|      13|Sammamish|  Bread|2017-01-07|       5|             null|
|       4|  Redmond|   Meat|2017-01-03|      40|             null|
|       7|  Redmond|  Bread|2017-01-05|       5|               40|
|       9|  Redmond| Cheese|2017-01-05|      15|                5|
|       1|  Seattle|Bananas|2017-01-01|       7|             null|
|       5|  Seattle|Potatoes|2017-01-04|      9|                7|
|       2|     Kent| Apples|2017-01-02|      20|             null|
|       3| Bellevue|Flowers|2017-01-02|      10|             null|
|       6| Bellevue|  Bread|2017-01-04|       5|               10|
|      18| Bellevue|  Bread|2017-01-11|       7|                5|
|      19| Bellevue|  Bread|2017-01-12|      54|                7|
+--------+---------+-------+----------+--------+-----------------+
```

## Pivot

```
[51]: spark.sql("select item, location from groceries_sql").groupBy("item").pivot("location").count().show()
```

```
+--------+--------+--------+----+-------+------+---------+-------+
|    item|Bellevue|Issaquah|Kent|Redmond|Renton|Sammamish|Seattle|
+--------+--------+--------+----+-------+------+---------+-------+
|Potatoes|    null|    null|null|   null|  null|     null|      1|
|  Cheese|    null|    null|null|      1|  null|     null|   null|
|    Meat|    null|       3|null|      1|  null|     null|   null|
|  Apples|    null|    null|   1|   null|  null|     null|   null|
|   Onion|    null|       3|null|   null|  null|     null|   null|
|   Bread|       5|    null|null|      1|     1|        1|   null|
| Flowers|       1|    null|null|   null|  null|     null|   null|
| Bananas|    null|    null|null|   null|  null|     null|      1|
|  Tomato|    null|       1|null|   null|  null|     null|   null|
+--------+--------+--------+----+-------+------+---------+-------+
```

```
[ ]:
```

# QUESTION 4

## Ways to handle Null Values : -

1. **Filtering** – With combination of filter() and isNotNull() function we can filter out the rows containing Null values.



2. **Dropping Null values** – The function drop() is used to remove rows containing Null values from your dataframes. This will create a new DataFrame without Null values.

**DF_without_Nulls = df.dropna()**

3. **Filling Null values -** fillna() is used to fill a value in place of null values present in particular column of a dataframe. We can use this to fill value in single column as well as all the columns containing nulls.

**Single Column – filled_df = df.fillna(replacement_value, subset = ['column name'])**

**All column – df.fillna(replacement_value)**

```
[33]: Lead_df = groceries_df.withColumn("Upcoming_quantity", lead("quantity").over(window_open_2))
```

```
[34]: Lead_df.show()
```

```
+--------+--------+-------+----------+--------+-----------------+
|order_id|location|   item|order_Date|quantity|Upcoming_quantity|
+--------+--------+-------+----------+--------+-----------------+
|       8|Issaquah|  Onion|2017-01-05|       4|                4|
|      10|Issaquah|  Onion|2017-01-06|       4|                4|
|      12|Issaquah|  Onion|2017-01-07|       4|                6|
|      14|Issaquah| Tomato|2017-01-07|       6|                3|
|      15|Issaquah|   Meat|2017-01-08|       3|                5|
|      16|Issaquah|   Meat|2017-01-09|       5|                6|
|      17|Issaquah|   Meat|2017-01-10|       6|             null|
|      13|Sammamish|  Bread|2017-01-07|      5|             null|
|       4| Redmond|   Meat|2017-01-03|      40|                5|
|       7| Redmond|  Bread|2017-01-05|       5|               15|
|       9| Redmond| Cheese|2017-01-05|      15|             null|
|       1| Seattle|Bananas|2017-01-01|       7|                9|
|       5| Seattle|Potatoes|2017-01-04|      9|             null|
|       2|    Kent| Apples|2017-01-02|      20|             null|
|       3|Bellevue|Flowers|2017-01-02|      10|                5|
|       6|Bellevue|  Bread|2017-01-04|       5|                7|
|      18|Bellevue|  Bread|2017-01-11|       7|               54|
|      19|Bellevue|  Bread|2017-01-12|      54|               34|
|      20|Bellevue|  Bread|2017-01-13|      34|               25|
|      21|Bellevue|  Bread|2017-01-14|      25|             null|
+--------+--------+-------+----------+--------+-----------------+
only showing top 20 rows
```

```
[42]: filled_df = Lead_df.fillna(0, subset = ["Upcoming_quantity"])
```

```
[43]: filled_df.show()
```

```
+--------+--------+-------+----------+--------+-----------------+
|order_id|location|   item|order_Date|quantity|Upcoming_quantity|
+--------+--------+-------+----------+--------+-----------------+
|       8|Issaquah|  Onion|2017-01-05|       4|                4|
|      10|Issaquah|  Onion|2017-01-06|       4|                4|
|      12|Issaquah|  Onion|2017-01-07|       4|                6|
|      14|Issaquah| Tomato|2017-01-07|       6|                3|
|      15|Issaquah|   Meat|2017-01-08|       3|                5|
|      16|Issaquah|   Meat|2017-01-09|       5|                6|
|      17|Issaquah|   Meat|2017-01-10|       6|                0|
|      13|Sammamish|  Bread|2017-01-07|      5|                0|
|       4| Redmond|   Meat|2017-01-03|      40|                5|
|       7| Redmond|  Bread|2017-01-05|       5|               15|
|       9| Redmond| Cheese|2017-01-05|      15|                0|
|       1| Seattle|Bananas|2017-01-01|       7|                9|
|       5| Seattle|Potatoes|2017-01-04|      9|                0|
|       2|    Kent| Apples|2017-01-02|      20|                0|
|       3|Bellevue|Flowers|2017-01-02|      10|                5|
|       6|Bellevue|  Bread|2017-01-04|       5|                7|
|      18|Bellevue|  Bread|2017-01-11|       7|               54|
|      19|Bellevue|  Bread|2017-01-12|      54|               34|
|      20|Bellevue|  Bread|2017-01-13|      34|               25|
|      21|Bellevue|  Bread|2017-01-14|      25|                0|
+--------+--------+-------+----------+--------+-----------------+
only showing top 20 rows
```

```
[36]: pivot_df = spark.sql("select item, location from groceries_sql").groupBy("item").pivot("location").count()
```

```
[37]: pivot_df.show()
```

```
+--------+--------+--------+----+-------+------+---------+-------+
|    item|Bellevue|Issaquah|Kent|Redmond|Renton|Sammamish|Seattle|
+--------+--------+--------+----+-------+------+---------+-------+
|Potatoes|    null|    null|null|   null|  null|     null|      1|
|  Cheese|    null|    null|null|      1|  null|     null|   null|
|    Meat|    null|       3|null|      1|  null|     null|   null|
|  Apples|    null|    null|   1|   null|  null|     null|   null|
|   Onion|    null|       3|null|   null|  null|     null|   null|
|   Bread|       5|    null|null|      1|     1|        1|   null|
| Flowers|       1|    null|null|   null|  null|     null|   null|
| Bananas|    null|    null|null|   null|  null|     null|      1|
|  Tomato|    null|       1|null|   null|  null|     null|   null|
+--------+--------+--------+----+-------+------+---------+-------+
```

```
[40]: filled_df = pivot_df.fillna(0)
```

```
[41]: filled_df.show()
```

```
+--------+--------+--------+----+-------+------+---------+-------+
|    item|Bellevue|Issaquah|Kent|Redmond|Renton|Sammamish|Seattle|
+--------+--------+--------+----+-------+------+---------+-------+
|Potatoes|       0|       0|   0|      0|     0|        0|      1|
|  Cheese|       0|       0|   0|      1|     0|        0|      0|
|    Meat|       0|       3|   0|      1|     0|        0|      0|
|  Apples|       0|       0|   1|      0|     0|        0|      0|
|   Onion|       0|       3|   0|      0|     0|        0|      0|
|   Bread|       5|       0|   0|      1|     1|        1|      0|
| Flowers|       1|       0|   0|      0|     0|        0|      0|
| Bananas|       0|       0|   0|      0|     0|        0|      1|
|  Tomato|       0|       1|   0|      0|     0|        0|      0|
+--------+--------+--------+----+-------+------+---------+-------+
```

4. **Nulls in Mathematical operations –** Use the na propery to handle nulls in mathematical operations. For example - null values can be replaced by a value when performing arithmetic operation.

**Math_df = df.na.fill(replacement_value)**

5. **Spark.sql query filter -** We can write custom spark.sql query to filter out null values using where column_name is not null.

```
[25]: Lead_df = groceries_df.withColumn("Upcoming_quantity", lead("quantity").over(window_open_2))

[26]: Lead_df.show()

+--------+--------+--------+----------+--------+-----------------+
|order_id|location|    item|order_Date|quantity|Upcoming_quantity|
+--------+--------+--------+----------+--------+-----------------+
|       8|Issaquah|   Onion|2017-01-05|       4|                4|
|      10|Issaquah|   Onion|2017-01-06|       4|                4|
|      12|Issaquah|   Onion|2017-01-07|       4|                6|
|      14|Issaquah|  Tomato|2017-01-07|       6|                3|
|      15|Issaquah|    Meat|2017-01-08|       3|                5|
|      16|Issaquah|    Meat|2017-01-09|       5|                6|
|      17|Issaquah|    Meat|2017-01-10|       6|             null|
|      13|Sammanish|   Bread|2017-01-07|       5|             null|
|       4| Redmond|    Meat|2017-01-03|      40|                5|
|       7| Redmond|   Bread|2017-01-05|       5|               15|
|       9| Redmond|  Cheese|2017-01-05|      15|             null|
|       1| Seattle| Bananas|2017-01-01|       7|                9|
|       5| Seattle|Potatoes|2017-01-04|       9|             null|
|       2|    Kent|  Apples|2017-01-02|      20|             null|
|       3|Bellevue| Flowers|2017-01-02|      10|                5|
|       6|Bellevue|   Bread|2017-01-04|       5|                7|
|      18|Bellevue|   Bread|2017-01-11|       7|               54|
|      19|Bellevue|   Bread|2017-01-12|      54|               34|
|      20|Bellevue|   Bread|2017-01-13|      34|               25|
|      21|Bellevue|   Bread|2017-01-14|      25|             null|
+--------+--------+--------+----------+--------+-----------------+
only showing top 20 rows
```

```
[46]: Lead_df.createOrReplaceTempView("Lead_sql")

[49]: spark.sql("select * from Lead_sql where Upcoming_quantity is not null").show()

+--------+--------+-------+----------+--------+-----------------+
|order_id|location|   item|order_Date|quantity|Upcoming_quantity|
+--------+--------+-------+----------+--------+-----------------+
|       8|Issaquah|  Onion|2017-01-05|       4|                4|
|      10|Issaquah|  Onion|2017-01-06|       4|                4|
|      12|Issaquah|  Onion|2017-01-07|       4|                6|
|      14|Issaquah| Tomato|2017-01-07|       6|                3|
|      15|Issaquah|   Meat|2017-01-08|       3|                5|
|      16|Issaquah|   Meat|2017-01-09|       5|                6|
|       4| Redmond|   Meat|2017-01-03|      40|                5|
|       7| Redmond|  Bread|2017-01-05|       5|               15|
|       1| Seattle|Bananas|2017-01-01|       7|                9|
|       3|Bellevue|Flowers|2017-01-02|      10|                5|
|       6|Bellevue|  Bread|2017-01-04|       5|                7|
|      18|Bellevue|  Bread|2017-01-11|       7|               54|
|      19|Bellevue|  Bread|2017-01-12|      54|               34|
|      20|Bellevue|  Bread|2017-01-13|      34|               25|
+--------+--------+-------+----------+--------+-----------------+
```