

INTRODUCTION TO HDFS:

HDFS is a storage system of Hadoop framework that is used to store files across a cluster of nodes by converting the file to multiple blocks of data and storing them with redundant copying mechanisms. Redundancy is introduced in the filing system because the nodes are usually hosted on hardware terminals which can go down at any moment, hence there must be a backup for all the data chunks that are created. HDFS works particularly well for very large files. The file sizes can be in the range of thousands of Gb to Tb's. HDFS works on the principle of allowing the computing power to go the data rather than getting the data to the computation machine when the size of the data is too large. Another major advantage of HDFS is that it can work on multiple machines supporting multiple software without any compatibility issues.

PRACTISING APACHE SPARK AND THE GIVEN QUERIES:

I used the project_utilities python file given in the mini project section of the canvas file section to read the data and to store it in the landing directory. I downloaded the data for the years 2004 and 2005 using the python script and then created a data lake directory for storing the data. I copied the file for the 2005 weather report and put it in the data lake hdfs directory on local host. Then I checked the number of records that were loaded to the data lake directory and decide to view the first few rows of the data to check for the column headings. View distinct event types helps us to view all the possible events that have occurred at least once in the data set. Then I saw how to subset the data only by choosing a few features out of all the available features. The next step that I took was to filter the data by the event type to find out the count and the occurrences of tornado.

Creating a Data Lake folder:

```
[js-170-211] ojaash ~-->hadoop fs -mkdir /data-lake
[js-170-211] ojaash ~-->hadoop fs -ls /
Found 1 items
drwxr-xr-x  - ojaash supergroup          0 2020-03-18 11:59 /data-lake
[js-170-211] ojaash ~-->
```

Transferring the file to the Data-Lake directory:

```
[js-170-211] ojaash ~-->hadoop fs -put /home/ojaash/landDir/StormEvents_details-ftp_v1.0_d2005_c20190920.csv.gz /data-lake
2020-03-18 12:11:50,216 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localhostTrusted = false, remoteHostTrusted = false
[js-170-211] ojaash ~-->hadoop fs -ls /data-lake
Found 1 items
-rw-r--r--  1 ojaash supergroup    7752031 2020-03-18 12:11 /data-lake/StormEvents_details-ftp_v1.0_d2005_c20190920.csv.gz
[js-170-211] ojaash ~-->
```

Checking the number of records loaded:

```
>>> df.count()
20/03/18 12:17:31 WARN Utils: Truncated the string representation of a plan since it was too large. This behavior can be adjusted by setting 'spark.debug.maxT
oStringFields' in SparkEnv.conf.
53976
```

Viewing the first 3 rows of the data:

[illegible]

Viewing distinct event types:

```
>>> df.select("EVENT_TYPE").distinct().show()
+-----+
|      EVENT TYPE|
+-----+
|    Winter Storm|
| Storm Surge/Tide|
|   Volcanic Ash|
| Marine High Wind|
|    Avalanche|
|   Dense Fog|
| Cold/Wind Chill|
|    High Surf|
| Tropical Storm|
|      Sleet|
|    High Wind|
|   Lightning|
| Coastal Flood|
|   Dust Devil|
| Winter Weather|
|    Wildfire|
|   Debris Flow|
|    Tornado|
| Frost/Freeze|
|   Dust Storm|
+-----+
only showing top 20 rows
```

Events by Year, Month Name and Locations

```
>>> df.select("YEAR", "MONTH_NAME", "STATE", "EVENT_TYPE").show()
```

YEAR	MONTH	NAME	STATE	EVENT	TYPE
2005	January	NEBRASKA	Heavy Snow		
2005	January	NEBRASKA	Heavy Snow		
2005	January	MONTANA	Heavy Snow		
2005	January	MONTANA	Heavy Snow		
2005	January	MONTANA	Heavy Snow		
2005	January	WISCONSIN	Winter Storm		
2005	January	VIRGINIA	High Wind		
2005	January	VIRGINIA	High Wind		
2005	January	VIRGINIA	High Wind		
2005	January	VIRGINIA	High Wind		
2005	January	WISCONSIN	Winter Storm		
2005	January	VIRGINIA	Ice Storm		
2005	March	ILLINOIS	Hail		
2005	January	VIRGINIA	Ice Storm		
2005	January	IOWA	High Wind		
2005	January	IOWA	High Wind		
2005	January	TEXAS	Hail		
2005	January	NEBRASKA	Heavy Snow		
2005	January	VIRGINIA	Ice Storm		
2005	January	VIRGINIA	Ice Storm		

only showing top 20 rows

Tornado events filtered

[illegible]

IMPLEMENTING A CUSTOM QUERY:

Question: Use the filter () function to count how many events of each type were reported by each source (SOURCE) in Georgia for your selected year. Write commands separately first and then chain them together.

Here I first decide to write the command to filter the data by the State name as Georgia. Then I decided to have a look at the number of records that exist for Georgia using the count () function.

```
>>> df_georgia_filter = df.filter(df.STATE == "GEORGIA")
>>> df_georgia_filter.count()
2141
>>> █
```

The next step was grouping the subset data for the Georgia state using the event types in order to get a count of the events that took place.

```
>>> df_georgia_group = df_georgia_filter.groupby(df.EVENT_TYPE).count()
>>> df_georgia_group.count()
21
█
```

Then I used the show function to display the data that I had created using the groupby command.

```
>>> df_georgia_group.show()
+-----+-----+
|          EVENT TYPE|count|
+-----+-----+
|      Winter Storm|   103|
| Tropical Storm   |   285|
|      High Surf   |     1|
|       Sleet      |     2|
|      High Wind   |    25|
|      Lightning   |    77|
| Winter Weather   |    38|
|      Wildfire    |     1|
|      Tornado     |    59|
|      Ice Storm   |    35|
| Funnel Cloud     |    17|
|      Flash Flood |    89|
| Freezing Fog     |    71|
|      Heavy Rain  |    18|
|       Heat       |    16|
|       Flood      |   122|
| Thunderstorm Wind|   286|
|       Hail       |   535|
|      Rip Current |     2|
|      Strong Wind |   139|
+-----+-----+
only showing top 20 rows
```

Finally, I combined the above syntax into a single statement and decided to print the data in ordering by the event types feature.

```
>>> df_georgia_agg = df.filter(df.STATE == "GEORGIA").groupby(df.MONTH_NAME, df.EVENT_TYPE).count().orderBy(df.EVENT_TYPE)
>>> df_georgia_agg.count()
87
>>> df_georgia_agg.show()
+-----+
|MONTH NAME| EVENT TYPE|count|
+-----+
| March| Flash Flood| 22|
| July| Flash Flood| 45|
| October| Flash Flood| 8|
| June| Flash Flood| 7|
| August| Flash Flood| 7|
| July| Flood| 57|
| June| Flood| 9|
| August| Flood| 8|
| March| Flood| 39|
| October| Flood| 5|
| February| Flood| 1|
| April| Flood| 3|
| December| Freezing Fog| 71|
| August| Funnel Cloud| 8|
| May| Funnel Cloud| 1|
| December| Funnel Cloud| 8|
| February| Hail| 88|
| November| Hail| 1|
| October| Hail| 9|
| April| Hail| 101|
+-----+
only showing top 20 rows
```

The implementation of the above queries was not that difficult. I could easily figure out the steps to be followed to subset the data and to apply the filters and groupby as well as orderBy commands. If we want we could have ordered the data to be displayed by the month name or the count as well. We can easily use the above code to study the data for different states as well as for different features and to get the proper conclusions.

CONCLUSION:

This assignment helps us to understand the different intricate features of HDFS which can be exploited to make the big data file storage, handling and analytics easier than any other framework. The Hadoop hdfs file systems can be accessed and queried very easily using the aggregate, groupby, orderby, filter, show, count and many more complex commands. The data can be subset and queried to get the results as desired as well as the data can be sent to any other machine for the visualization application and storage.