

Continuous Internal Assessment Test - II

Name : Shyam Sundar S
Roll No : 24MX347

Step 1 :

Run the command: hdfs dfs -ls /data/input

And create the sales.csv and weather.csv files.

Step 2:

Load the local csv files into Ubuntu hdfs storage

Screenshot of the output :

```
*****  
devcontainers@SHYAM:~$ start-dfs.sh  
Starting namenodes on [localhost]  
Starting datanodes  
localhost: datanode is running as process 12611. Stop it first and ensure /tmp/hadoop-devcontainers-datanode.pid file is empty before retry.  
Starting secondarynamenodes [SHYAM]  
SHYAM: secondarynamenode is running as process 12784. Stop it first and ensure /tmp/hadoop-devcontainers-secondarynamenode.pid file is empty before retry.  
devcontainers@SHYAM:~$ start-yarn.sh  
Starting resourcemanager  
resourcemanager is running as process 13087. Stop it first and ensure /tmp/hadoop-devcontainers-resourcemanager.pid file is empty before retry.  
Starting nodemanagers  
localhost: nodemanager is running as process 13224. Stop it first and ensure /tmp/hadoop-devcontainers-nodemanager.pid file is empty before retry.  
devcontainers@SHYAM:~$ hdfs dfs -mkdir data/input/  
devcontainers@SHYAM:~$ hdfs dfs -put sales.csv /data/input/  
devcontainers@SHYAM:~$ hdfs dfs -put weather.csv /data/input/  
devcontainers@SHYAM:~$ hdfs dfs -ls /data/input/  
Found 2 items  
-rw-r--r-- 1 devcontainers supergroup 197 2025-10-13 04:17 /data/input/sales.csv  
-rw-r--r-- 1 devcontainers supergroup 178 2025-10-13 04:17 /data/input/weather.csv  
devcontainers@SHYAM:~$ nano analysis.pig  
devcontainers@SHYAM:~$ pig -x mapreduce  
2025-10-13 04:18:39.695 INFO ExecTypeProvider: Trying ExecType : LOCAL  
2025-10-13 04:18:39.697 INFO ExecTypeProvider: Trying ExecType : MAPREDUCE  
2025-10-13 04:18:39.697 INFO pig ExecTypeProvider: Picked MAPREDUCE as the ExecType  
2025-10-13 04:18:39.758 [main] INFO org.apache.pig.Main - Apache Pig version 0.17.0 (r1797386) compiled Jun 02 2017, 15:41:58  
2025-10-13 04:18:39.758 [main] INFO org.apache.pig.Main - Logging error messages to: /home/devcontainers/pig_1760329110753.log  
2025-10-13 04:18:39.999 [main] INFO org.apache.pig.impl.util.Utils - Default bootstrap file /home/devcontainers/.pigbootup found  
2025-10-13 04:18:31.060 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address  
2025-10-13 04:18:31.060 [main] INFO org.apache.hadoop.mapred.lib.HadoopConfiguration - Connecting to hadoop file system at: hdfs://localhost:9000  
2025-10-13 04:18:31.538 [main] INFO org.apache.pig.PigServer - Pig Script ID for the session: PIG-analysis.pig-acfed3e-6c18-4c2c-95b1-acbf9e771d95  
2025-10-13 04:18:31.538 [main] WARN org.apache.pig.PigServer - ATS is disabled since yarn.timeline-service.enabled set to false
```

Pig Analysis Generation using pig Latin script Screenshot:

```
HadoopVersion PigVersion UserId StartedAt FinishedAt Features  
3.3.6 0.17.0 devcontainers 2025-10-13 04:33:27 2025-10-13 04:33:52 HASH_JOIN,GROUP_BY,ORDER_BY,FILTER,LIMIT  
  
Success:  
Job Stats (time in seconds):  
Job: Maps 0 Reduces 0 MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReduceTime Alias Feature Outputs  
job_local1030950535_0001 1 n/a n/a n/a n/a n/a n/a n/a sorted_sales SAMPLER HASH_JOIN,MULTI_QUERY hdfs  
job_local1268590773_0002 2 1 n/a n/a n/a n/a n/a n/a sales_weather JOIN,sunny_day_sales,weather MULTI_QUERY,COMBINER hdfs  
://data/output/sunny_day_sales,  
job_local13574406 0001 1 n/a n/a n/a n/a n/a n/a n/a high_value_sales,sales,sales_by_product,total_sales MULTI_QUERY,COMBINER hdfs  
://data/output/high_value_sales,hdfs://data/output/total_sales,  
job_local1571309557_0004 1 1 n/a n/a n/a n/a n/a n/a n/a sorted_sales ORDER_BY,COMBINER hdfs  
job_local832496057_0005 1 1 n/a n/a n/a n/a n/a n/a avg_sales_by_condition,grouped_by_condition GROUP_BY,COMBINER hdfs://data  
/output/avg_sales_by_condition,  
job_local856747376_0006 1 1 n/a n/a n/a n/a n/a n/a n/a sorted_sales hdfs://data/output/top_3_sales,  
  
Input(s):  
Successfully read 7 records (11511958 bytes) from: "hdfs://data/input/sales.csv"  
Successfully read 7 records from: "hdfs://data/input/weather.csv"  
  
Output(s):  
Successfully stored 2 records (63 bytes) in: "hdfs://data/output/high_value_sales"  
Successfully stored 3 records (35 bytes) in: "hdfs://data/output/total_sales"  
Successfully stored 3 records (49971963 bytes) in: "hdfs://data/output/sunny_day_sales"  
Successfully stored 3 records (167 bytes) in: "hdfs://data/output/sunny_day_sales"  
Successfully stored 4 records (57560267 bytes) in: "hdfs://data/output/avg_sales_by_condition"  
  
Counters:  
Total records written : 15  
Total bytes written : 126632479  
Spilled Local Memory (bytes) : 11  
Spilled Local Memory (bytes) until count : 0  
Total bytes proactively spilled: 0  
Total records proactively spilled: 0  
  
Job DAG:  
job_local13574406_0001 -> job_local1030950535_0003,job_local1268590773_0002,  
job_local1030950535_0003 -> job_local1571309557_0004,  
job_local1571309557_0004 -> job_local856747376_0006,  
job_local856747376_0006 -> job_local832496057_0005,  
job_local832496057_0005 ->  
  
2025-10-13 04:33:52,213 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!  
2025-10-13 04:33:52,214 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!  
2025-10-13 04:33:52,216 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!  
2025-10-13 04:33:52,220 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!  
2025-10-13 04:33:52,222 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
```

Step 3:

Queries of the Script :

What is the total sales amount for each product?

Which sales transactions had a value greater than \$200?

Which sales occurred on days when the weather was sunny?

What is the average sales amount for each weather condition?

What are the top 3 highest sales transactions by amount?

Five Queries Output Screenshot:

Run the command: hdfs dfs -cat /data/output/total_sales/part-r-00000

```
devcontainers@SHYAM:~$ hdfs dfs -cat /data/output/total_sales/part-r-00000
P001, 400.25
P002, 325.5
P003, 550.75
devcontainers@SHYAM:~$ hdfs dfs -cat /data/output/total_sales/part-00000
cat: '/data/output/total_sales/part-00000': No such file or directory
devcontainers@SHYAM:~$ hdfs dfs -getmerge /data/output/* all_outputs.txt
devcontainers@SHYAM:~$ cat all_outputs.txt
Rain,150.0
Clear,250.75
Sunny,208.5
Cloudy,125.125
P003, 20251011,300.0,Tokyo
P003, 20251012,250.75,Tokyo
P003, 20251011,300.0,Tokyo,Tokyo,20251011,25,Sunny
P002, 20251010,150.5,New York,New York,20251010,22,Sunny
P002, 20251012,175.0,New York,New York,20251012,24,Sunny
P003, 20251011,300.0,Tokyo
P003, 20251012,250.75,Tokyo
P001, 20251010,200.0,London
P001, 400.25
P002, 325.5
P003, 550.75
devcontainers@SHYAM:~$ |
```