

Snapshots of steps :

Create Cluster - Quick Options [Go to advanced options](#)

General Configuration

Cluster name

☒ Logging ⓘ
S3 folder 📁

Launch mode ☒ Cluster ⓘ ☐ Step execution ⓘ

Software configuration

Release ⓘ

Applications ☒ Core Hadoop: Hadoop 3.2.1 with Hive 3.1.2, Hue 4.9.0, Pig 0.17.0 and Tez 0.9.2
☐ HBase: HBase 2.2.6 with Hadoop 3.2.1, Hive 3.1.2,

aws

Services

Search

[Alt+S]

Amazon EMR

EMR Studio

EMR Serverless New

EMR on EC2

Clusters

Notebooks

Git repositories

Security configurations

Block public access

VPC subnets

Events

EMR on EKS

Virtual clusters

Help

Clone Terminate AWS CLI export

Cluster: project_shreeya_desai Starting

Summary Application user interfaces Monitoring Hardware Configurations Events Steps Bootstrap actions

Summary

ID: j-27OAZUD4I73BX
Creation date: 2023-04-03 22:55 (UTC-4)
Elapsed time: 0 seconds
After last step completes: Cluster waits
Termination protection: Off [Change](#)
Tags: -- [View All / Edit](#)
Master public DNS: --

Configuration details

Release label: emr-6.3.1
Hadoop distribution: Amazon 3.2.1
Applications: Hive 3.1.2, Hue 4.9.0, Pig 0.17.0, Tez 0.9.2
Log URI: s3://aws-logs-390943627784-us-east-2/elasticmapreduce/ 📁
EMRFS consistent view: Disabled

CloudShell Feedback Language © 2023, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

aws

Services

Search

[Alt+S]

Upload succeeded

View details below.

Summary

Destination

s3://aws-logs-390943627784-us-east-2/project1/

Succeeded

✔ 1 file, 515.3 MB (100.00%)

Failed

✖ 0 files, 0 B (0%)

Files and folders

Configuration

Files and folders (1 Total, 515.3 MB)

< 1 >

Name	Folder	Type	Size	Status	Error
1997.csv	-	text/csv	515.3 MB	✔ Succeeded	-

CloudShell Feedback Language © 2023, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences


```
hive> SELECT Origin AS Airport, SUM((ArrDelay + DepDelay) / 60.0) AS DelayHours FROM AIRLINESPROJECT_SHREEYA WHERE Origin IS NOT NULL AND ArrDelay IS NOT NULL AND DepDelay IS NOT NULL GROUP BY Origin ORDER BY DelayHours DESC LIMIT 3;
Query ID = hadoop_20230404034219_01bd0086-e276-462a-9f6e-910df977db84
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1680577291903_0006)
```

VERTICES	MODE	STATUS	TOTAL	COMPLETED	RUNNING	PENDING	FAILED	KILLED
Map 1	container	SUCCEEDED	1	1	0	0	0	0
Reducer 2	container	SUCCEEDED	2	2	0	0	0	0
Reducer 3	container	SUCCEEDED	1	1	0	0	0	0

```
VERTICES: 03/03 [=====]>>] 100% ELAPSED TIME: 12.61 s
-----
OK
ATL      88807.616761
ORD      84955.366468
DFW      67215.133298
Time taken: 15.293 seconds, Fetched: 3 row(s)
hive>
```

```
hive>
> SELECT UniqueCarrier AS Carrier, SUM(IF(ArrDelay > DepDelay, ArrDelay, DepDelay)) / 60.0 AS DelayHours FROM AIRLINESPROJECT_SHREEYA WHERE ArrDelay IS NOT NULL AND DepDelay IS NOT NULL GROUP BY UniqueCarrier ORDER BY DelayHours DESC LIMIT 3;
Query ID = hadoop_20230404034314_9cc7c08b-e9c4-4d1c-bb4e-ec3f3b61379c
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1680577291903_0006)
```

VERTICES	MODE	STATUS	TOTAL	COMPLETED	RUNNING	PENDING	FAILED	KILLED
Map 1	container	SUCCEEDED	1	1	0	0	0	0
Reducer 2	container	SUCCEEDED	2	2	0	0	0	0
Reducer 3	container	SUCCEEDED	1	1	0	0	0	0

```
VERTICES: 03/03 [=====]>>] 100% ELAPSED TIME: 7.86 s
-----
OK
DL      205310.200000
UA      162923.200000
WN      133664.450000
Time taken: 8.986 seconds, Fetched: 3 row(s)
hive>
```

```
hive> SELECT 'Arrival Delay' AS DelayType, SUM(ArrDelay) AS TotalDelayMinutes FROM AIRLINESPROJECT_SHREEYA WHERE ArrDelay IS NOT NULL UNION ALL SELECT 'Departure Delay' AS DelayType, SUM(DepDelay) AS TotalDelayMinutes FROM AIRLINESPROJECT_SHREEYA WHERE DepDelay IS NOT NULL ORDER BY TotalDelayMinutes DESC LIMIT 1;
Query ID = hadoop_20230404042830_1ffbd209-bd54-4c6e-839c-5d9ecfe76fb1
Total jobs = 1
Launching Job 1 out of 1
Tsz session was closed. Reopening...
Session re-established.
Session re-established.
Status: Running (Executing on YARN cluster with App id application_1680577291903_0013)
```

VERTICES	MODE	STATUS	TOTAL	COMPLETED	RUNNING	PENDING	FAILED	KILLED
Map 1	container	SUCCEEDED	1	1	0	0	0	0
Reducer 2	container	SUCCEEDED	1	1	0	0	0	0
Reducer 4	container	SUCCEEDED	1	1	0	0	0	0
Map 5	container	SUCCEEDED	1	1	0	0	0	0
Reducer 6	container	SUCCEEDED	1	1	0	0	0	0

```
VERTICES: 05/05 [=====]>>] 100% ELAPSED TIME: 7.51 s
-----
OK
Departure Delay 43764457
Time taken: 13.081 seconds, Fetched: 1 row(s)
hive>
```

Commands used:

1. `wget (s3 path)`
2. `hadoop fs -mkdir /shreeya`
3. `hadoop fs -mkdir /shreeya/csv`
4. `hadoop fs -put 1997.csv /shreeya/csv`
5. `hive`
6. `CREATE EXTERNAL TABLE IF NOT EXISTS AIRLINESPROJECT_SHREEYA (Year INT, Month INT, DayofMonth INT, DayOfWeek INT, DepTime INT, CRSDepTime INT, ArrTime INT, CRSArrTime INT, UniqueCarrier STRING, FlightNum INT, TailNum STRING, ActualElapsedTime INT, CRSElapsedTime INT, AirTime INT, ArrDelay INT, DepDelay INT, Origin STRING, Dest STRING, Distance INT, TaxiIn INT, TaxiOut INT, Cancelled INT, CancellationCode STRING, Diverted INT, CarrierDelay INT, WeatherDelay INT, NASDelay INT, SecurityDelay INT, LateAircraftDelay INT) ROW FORMAT DELIMITED FIELDS TERMINATED BY "," tblproperties("skip.header.line.count"="1");`
7. `load data inpath '/shreeya/csv/1997.csv' into table AIRLINESPROJECT_SHREEYA;`
8. `SELECT * FROM AIRLINESPROJECT_SHREEYA LIMIT 15;`

9. **QUERY TO DETERMINE THE THREE AIRPORTS WITH THE HIGHEST DELAY TIME (IN HOURS):**

```
SELECT Origin AS Airport, SUM((ArrDelay + DepDelay) / 60.0) AS DelayHours FROM  
AIRLINESPROJECT_SHREEYA WHERE Origin IS NOT NULL AND ArrDelay IS NOT NULL AND  
DepDelay IS NOT NULL GROUP BY Origin ORDER BY DelayHours DESC LIMIT 3;
```

10. **QUERY TO determine the three carriers with the highest delay time (in hours):**

```
SELECT UniqueCarrier AS Carrier, SUM((ArrDelay + DepDelay) / 60.0) AS DelayHours FROM  
AIRLINESPROJECT_SHREEYA WHERE UniqueCarrier IS NOT NULL AND ArrDelay IS NOT NULL  
AND DepDelay IS NOT NULL GROUP BY UniqueCarrier ORDER BY DelayHours DESC LIMIT 3;
```

11. **QUERY TO determine overall which type of delay (arrivals or departures) is the largest for airports:**

```
SELECT 'Arrival Delay' AS DelayType, SUM(ArrDelay) AS TotalDelayMinutes FROM  
AIRLINESPROJECT_SHREEYA WHERE ArrDelay IS NOT NULL UNION ALL SELECT 'Departure  
Delay' AS DelayType, SUM(DepDelay) AS TotalDelayMinutes FROM  
AIRLINESPROJECT_SHREEYA WHERE DepDelay IS NOT NULL ORDER BY  
TotalDelayMinutes DESC LIMIT 1;
```