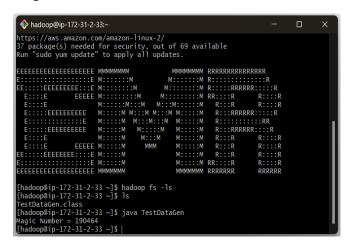
CSP554 Big Data Technologies

Assignment 4

ID: A20522183

1. java TestDataGen

Magic Number = 190464



2. placed foodata and foodratings

```
[hadoop@ip-1/2-31-2-33 ~]$ java TestDataGen
Magic Number = 190464
[hadoop@ip-172-31-2-33 ~]$ ls
foodplaces190464.txt foodratings190464.txt TestDataGen.class
[hadoop@ip-172-31-2-33 ~]$
```

3. View data in food places

```
E EEEEE
EEEEEEEE:::E
```

Create a database on HIVE

Create Database mydb.

```
Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j2.
properties Async: false
hive> CREATE DATABASE MyDb;
Time taken: 1.967 seconds
```

```
Exercise 1.

Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j2.

properties Async: false
hive> CREATE DATABASE MyDb;

OK
Time taken: 1.967 seconds
hive>
```

Created a database on HIVE

```
java:62)
    at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMethodAcces
sorImpl.java:43)
    at java.lang.reflect.Method.invoke(Method.java:498)
    at org.apache.hadoop.util.RunJar.main(RunJar.java:244)
    at org.apache.hadoop.util.RunJar.main(RunJar.java:158)

FAILED: ParseException line 11:10 mismatched input '/' expecting StringLiteral n ear 'LOCATION' in table location specification
hive> Create External table if not exists MyDb.foodratings(
    > name String Comment 'Food Critic Name',
    > food1 INT Comment 'Review Rating 1',
    > food2 INT Comment 'Review Rating 2',
    > food3 INT Comment 'Review Rating 3',
    > food4 INT Comment 'Review Rating 4',
    > id INT Comment 'Restaurant Id(FK)')
    > Comment 'Ratings Data'
    > ROW FORMAT DELIMITED FIELDS TERMINATED BY ','
    > STORED AS TEXTFILE
    > LOCATION '/home/hadoop/foodratings190464.txt';
OK
Time taken: 0.267 seconds
hive>
```

Describe formatted mydb.foodratings;

```
hadoop@ip-172-31-2-33:~
                                                                            X
ime taken: 0.267 seconds
nive> DESCRIBE FORMATTED MyDb.foodratings;
 col_name
                         data_type
                                                   comment
                         string
                                                   Food Critic Name
name
ood1
                                                   Review Rating 1
                         int
ood2
                                                   Review Rating 2
                         int
Food3
                         int
                                                   Review Rating 3
ood4
                         int
                                                   Review Rating 4
i d
                         int
                                                   Restaurant Id(FK)
 Detailed Table Information
Database:
                         mydb
Owner:
                         hadoop
reateTime:
                         Thu Sep 28 16:46:07 UTC 2023
_astAccessTime:
Retention:
                         UNKNOWN
ocation:
                         hdfs://ip-172-31-2-33.ec2.internal:8020/home/hadoop/food
atings190464.txt
Table Type:
                         EXTERNAL_TABLE
Table Parameters:
       EXTERNAL
                                 TRUE
       comment
                                 Ratings Data
```

```
hadoop@ip-172-31-2-33:~
                                                                          ×
astAccessTime:
                        UNKNOWN
Retention:
                        hdfs://ip-172-31-2-33.ec2.internal:8020/home/hadoop/food
_ocation:
ratings190464.txt
Table Type:
                        EXTERNAL_TABLE
able Parameters:
        EXTERNAL
                                TRUE
        comment
                                 Ratings Data
        transient_lastDdlTime
                                1695919567
# Storage Information
SerDe Library:
                        org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe
InputFormat:
                        org.apache.hadoop.mapred.TextInputFormat
OutputFormat:
                        org.apache.hadoop.hive.ql.io.HiveIgnoreKeyTextOutputForm
at
Compressed:
                        No
Num Buckets:
                        -1
                        Bucket Columns:
Sort Columns:
Storage Desc Params:
        field.delim
        serialization.format
Time taken: 0.127 seconds, Fetched: 33 row(s)
hive>
```

```
hadoop@ip-172-31-2-33:~
                                                                         X
Database:
                        mydb
Owner:
                        hadoop
                        Thu Sep 28 16:51:39 UTC 2023
reateTime:
.astAccessTime:
                        UNKNOWN
Retention:
ocation:
                        hdfs://ip-172-31-2-33.ec2.internal:8020/home/hadoop/food
places190464.txt
Table Type:
                        EXTERNAL_TABLE
Table Parameters:
       EXTERNAL
                                TRUE
                                Restaurant Details
        comment
        transient_lastDdlTime
                                1695919899
 Storage Information
SerDe Library:
                        org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe
InputFormat:
                        org.apache.hadoop.mapred.TextInputFormat
OutputFormat:
                        org.apache.hadoop.hive.gl.io.HiveIgnoreKeyTextOutputForm
at
                        No
Compressed:
Num Buckets:
                        Bucket Columns:
Sort Columns:
Storage Desc Params:
        field.delim
        serialization.format
Time taken: 0.058 seconds, Fetched: 29 row(s)
hive>
```

```
hadoop@ip-172-31-2-33:~
                                                                                                       X
Num Buckets:
                            -1
[]
[]
Bucket Columns:
Sort Columns:
Storage Desc Params:
field.delim
         serialization.format
Time taken: 0.058 seconds, Fetched: 29 row(s)
hive> select name, min(food3) as MIN, max(food3) as MAX, avg(food3) as AVG from
MyDb.foodratings;
FAILED: SemanticException [Error 10025]: Line 1:7 Expression not in GROUP BY key
hive> select "food3" as Column_name, min(food3) as MIN, max(food3) as MAX, avg(food3) as AVG from MyDb.foodratings;
Query ID = hadoop_20230928165721_d74a90f9-538d-4b4d-a67b-e4d1bd9a0269
Total jobs = 1
Launching Job 1 out of 1
Tez session was closed. Reopening...
Session re-established.
Status: Running (Executing on YARN cluster with App id application_1695917019180
_0002)
Map 1: -/-
-Map 1: -/-
Map 1: -/-
                   Reducer 2: 0/1
                  Reducer 2: 0(+1)/1
Reducer 2: 1/1
oĸ
food3 NULL
                   NULL
                          NULL
Time taken: 9.586 seconds, Fetched: 1 row(s)
hive>
```

Select name, min(food1) as min, max(food1) as max, avg(food1) as avg from mydb.foodratings group by name.

```
hive> SELECT name, min(food1) AS MIN, max(food1) as MAX, AVG(food1) as AVG from MyDb.foodratings Group ny name;
FAILED: ParseException line 1:97 missing BY at 'ny' near '<EOF>'
line 1:100 extraneous input 'name' expecting EOF near '<EOF>'
hive> SELECT name, min(food1) AS MIN, max(food1) as MAX, AVG(food1) as AVG from MyDb.foodratings Group by name;
Query ID = hadoop_20221001083424_93670229-6565-4bd4-b9e9-400b38a1ed61
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1664609429932_0002)
       VERTICES MODE STATUS TOTAL COMPLETED RUNNING PENDING FAILED KILLED
Map 1 ...... container SUCCEEDED 1
Reducer 2 ..... container SUCCEEDED 2
                                                                           0
                                                                                           0
                                                                  a
                                                                                   0
 ······
 'ERTICES: 02/02 [==============>>] 100% ELAPSED TIME: 4.61 s
Jill
                       24.737373737373737
                        24.5555555555555
Joe
                       25.646464646464647
Joy
               50
Mel
                       24.989583333333333
               50
                       24.653658536585365
Sam
Time taken: 5.285 seconds, Fetched: 5 row(s)
```

```
hive> DESCRIBE FORMATTED MyDB.foodratingspart;
 col_name
                          data_type
                                                    comment
food1
                          int
                                                    Review rating 1
                                                    Review rating 2
food2
                          int
food3
                          int
                                                    Review rating
food4
                                                    Review rating 4
                          int
                          int
                                                    Restaurant ID (FK)
 Partition Information
 col_name
                          data_type
                                                    comment
                                                    Name of food Critic
                          string
# Detailed Table Information
Database:
Owner:
                          hadoop
reateTime:
                          Thu Sep 28 17:10:17 UTC 2023
LastAccessTime:
                          UNKNOWN
Retention:
Location:
                          hdfs://ip-172-31-2-33.ec2.internal:8020/user/hive/warehouse/mydb.db/foodratings
part
                          EXTERNAL_TABLE
Table Type:
Table Parameters:
        COLUMN_STATS_ACCURATE
                                   {\"BASIC_STATS\":\"true\"}
        EXTERNAL
                                   TRUE
        comment
numFiles
                                  Rating data
        numPartitions
        numRows
        rawDataSize
                                   0
        totalSize
        transient_lastDdlTime
                                   1695921017
 Storage Information
SerDe Library:
                          org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe
                          org.apache.hadoop.mapred.TextInputFormat org.apache.hadoop.hive.ql.io.HiveIgnoreKeyTextOutputFormat
InputFormat:
OutputFormat:
Compressed:
                          No
                          -1
[]
Num Buckets:
Bucket Columns:
Sort Columns:
Storage Desc Params:
        field.delim
serialization.format
ime taken: 0.049 seconds, Fetched: 43 row(s)
```

5

Answer: Partition on critic names would help querying faster and easier than partitioning on data of number of places which are large.

```
Time taken: 0.049 seconds, Fetched: 43 row(s)
hive> SET hive.exec.dynamic.partition.mode = non-strict;
hive> SET hive.exec.dynamic.partition;
hive.exec.dynamic.partition=true
hive> SET hive.exec.dynamic.partition.mode;
hive.exec.dynamic.partition.mode=non-strict
hive> |
```

Set hive.exec.dynamic.partiton.mode = non-strict;

```
INSERT OVERWRITE TABLE MyDb.foodratingspart
      PARTITION (name)
> SELECT food1,food2,food3,food4,id,name FROM MyDb.foodratings;
Query ID = hadoop_20230928171547_fcf547e0-b3df-43e1-bf85-da64c061a8c4
Fotal jobs = 1
aunching Job 1 out of 1
Tez session was closed. Reopening...
Session re-established.
        VERTICES
                       MODE
                                       STATUS TOTAL COMPLETED RUNNING PENDING FAILED KILLED
 ap 1
                   container
                                    SUCCEEDED
 ERTICES: 00/01 [>>-
                                                ----] 0%
                                                               ELAPSED TIME: 0.47 s
oading data to table mydb.foodratingspart partition (name=null)
          Time taken to load dynamic partitions: 0.009 seconds
          Time taken for adding to write entity: 0.0 seconds
Time taken: 5.447 seconds
nive>
```

INSERT OVERWRITE TABLE Mydb.foodratingspart.

SELECT MIN(FOOD1) AS MIN, MAX(FOOD1) AS MAX, AVG(FOOD1) AS AVG FROM MYDB.FOODRATINGSPART WHERE NAME = "JILL" OR NAME="MEL"

SELECT FOODP.PLACE, AVG(FOODR,FOOD4) AS AVG FROM MYDB.FOODP.ID =FOODP.ID GROUP BY FOODP.PLACE = 'SOUP BOWL' AND FOODP.ID =FOODR.ID GROUP BY

```
ive> select foodp.place, avg(foodr.food4) as AVG from mydb.foodplaces foodp, Mydb.foodratings foodr where foodp.place
Soup Bowl" and foodp.id =foodr.id group by foodp.place;
Query ID = hadoop_20221001085612_6756f362-678f-493b-8035-99bf080a8256
otal jobs = 1
aunching Job 1 out of 1
tatus: Running (Executing on YARN cluster with App id application_1664609429932_0003)
                  MODE STATUS TOTAL COMPLETED RUNNING PENDING FAILED KILLED
       .....
                           SUCCEEDED
                                                            0
                                                                            0
                                                                                   0
ap 1 ..... container
    ..... container
                           SUCCEEDED
                                                            0
                                                                            0
                                                                                   0
ducer 3 ..... container SUCCEEDED
             25.21578947368421
oup Bowl
ime taken: 8.056 seconds, Fetched: 1 row(s)
```

8.

a) Row format is useful when user has to access data w.r.t row values and need to access many rows at a time. This is for optimal data reading and writing.

Column format is useful when computation is focused on specific column without the need to search row values. This is for optimal data computation.

- b) The capability to divide a file into independent smaller components is termed "splitability." When it comes to columnar file formats, splitability is achievable when query processing is concentrated on a specific column. This allows the data to be partitioned based on columns, enhancing computational efficiency.
- c) Storing similar data types adjacent to each other enables more efficient compression compared to organizing them in a row-wise manner.
- d) Parquet is used in Hadoop analytical database like (Implaca). It is specially used in analysing huge dataset with multiple columns for computations.