# AADHAR

## CHECKPOINT 1

Load the data into HDFS, Hive Managed table, Hive External table and Spark DataFrame.

LOADING DATA IN HDFS

hdfs dfs -mkdir aadhar

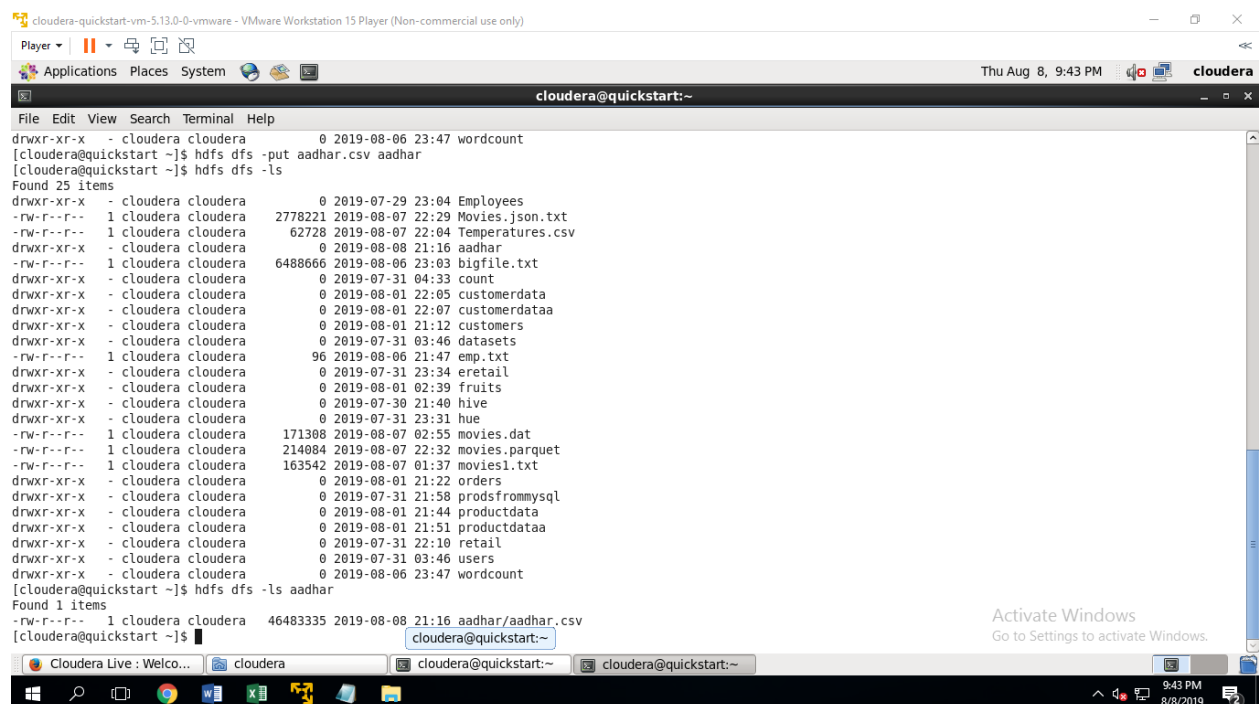hdfs dfs -put aadhar.csv aadhar

hdfs dfs -ls

hdfs dfs -ls aadhar

Found 1 items

-rw-r--r--   1 cloudera cloudera   46483335 2019-08-08 21:16 aadhar/aadhar.csv



for loading data in Hive managed table

start Hive

hive> create database if not exists aadhar;

OK

Time taken: 1.435 seconds

hive> use aadhar;

OK

Time taken: 0.082 seconds

create table if not exists aadhar_in(registrar String,enrolment_Agency String,State String,district String,sub_District String,pin_code String,gender String,age int,aadhaarGenerated int,enrolmentRejected int,residentsProvidingEmail int,residentsProvidingMobileNumber int)

    > row format delimited fields terminated by ','

    > stored as textfile

    > TBLPROPERTIES('serialization.null.format'='',
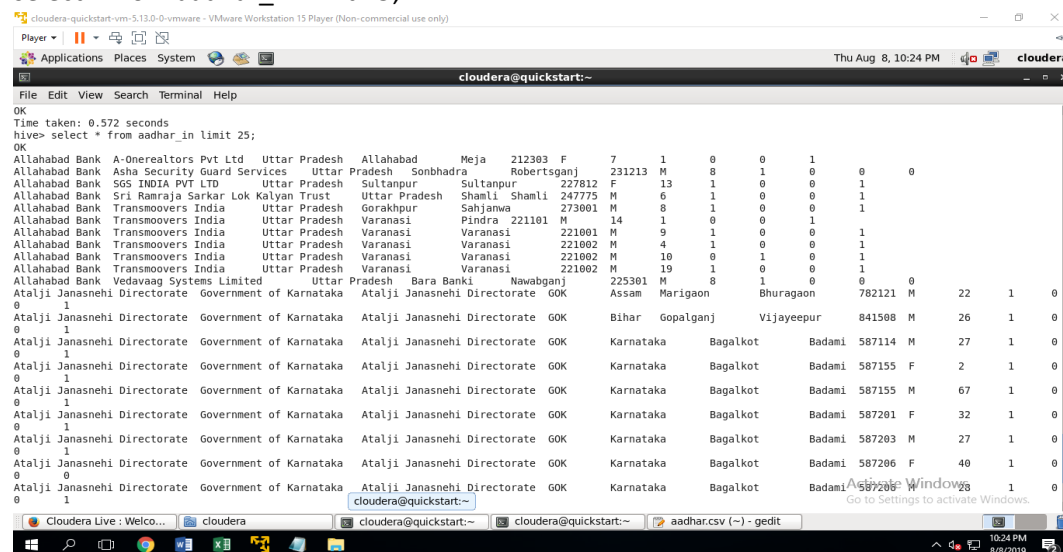
    > 'skip.header.line.count'='1');

OK

Time taken: 0.071 seconds

load data inpath '/user/cloudera/aadhar/aadhar.csv' overwrite into table aadhar_in;

select * from aadhar_in limit 25;



# For creation of hive external table

hive> create  external table if not exists aadhar_ex(registrar String,enrolment_Agency String,State String,district String,sub_District String,pin_code String,gender String,age int,aadhaarGenerated int,enrolmentRejected int,residentsProvidingEmail int,residentsProvidingMobileNumber int)

   > row format delimited fields terminated by ','

   > stored as textfile

   > location '/user/cloudera/aadhar/'

   > TBLPROPERTIES('serialization.null.format'='',

   > 'skip.header.line.count'='1');

OK

Time taken: 0.056 second



## Creation of the table into DF

val aadharrdd=sc.textFile("/user/cloudera/aadhar/aadhar.csv")

val header=aadharrdd.first();

header: String = Registrar,Enrolment Agency,State,District,Sub District,Pin Code,Gender,Age,Aadhaar generated,Enrolment Rejected,Residents providing email,Residents providing mobile number

val firstaadharrdd=aadharrdd.filter(row=>row!=header);

val aadharDF=firstaadharrdd.map(_.split(",")).map{case Array(a,b,c,d,e,f,g,h,i,j,k,l) => (a,b,c,d,e,f,g,h.toInt,i.toInt,j.toInt,k.toInt,l.toInt)}.toDF("registrar","enrollmentAgency","state","district","subDistrict","pinCode","gender","age","aadharGenerated","enrolmentRejected","residentsProvidingEmail","residentsProvidingMobileNumber");

aadharDF: org.apache.spark.sql.DataFrame = [registrar: string, enrollmentAgency: string, state: string, district: string, subDistrict: string, pinCode: string, gender: string, age: int, aadharGenerated: int, enrolmentRejected: int, residentsProvidingEmail: int, residentsProvidingMobileNumber: int]

aadharDF.show(25);



# CHECKPOINT 2

## 2. Describe schema

scala> aadharDF.printSchema

root

 |-- registrar: string (nullable = true)

 |-- enrollmentAgency: string (nullable = true)

 |-- state: string (nullable = true)

 |-- district: string (nullable = true)

 |-- subDistrict: string (nullable = true)

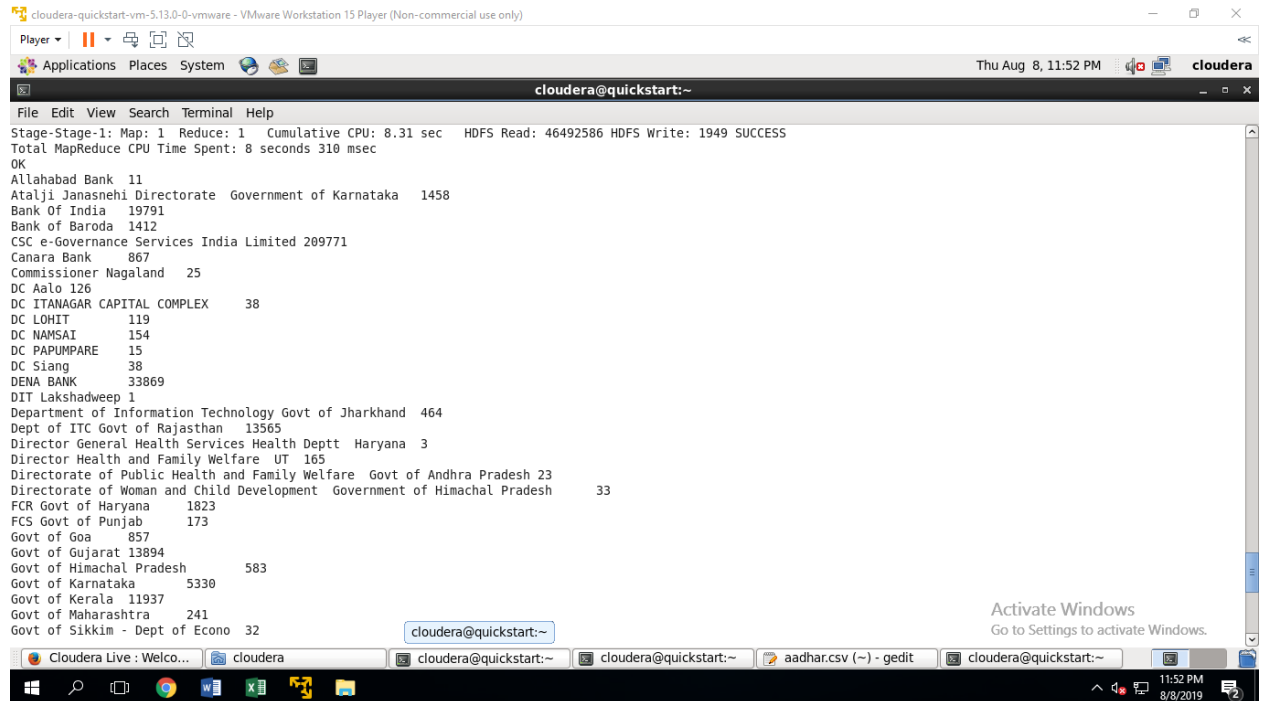|-- pinCode: string (nullable = true)

|-- gender: string (nullable = true)

|-- age: integer (nullable = false)

|-- aadharGenerated: integer (nullable = false)

|-- enrolmentRejected: integer (nullable = false)

|-- residentsProvidingEmail: integer (nullable = false)

|-- residentsProvidingMobileNumber: integer (nullable = false)



# 3. Find the count and names of registrars in the table

select registrar,count(*) from aadhar_ex groupby registrar;


OK

Allahabad Bank 11

Atalji Janasnehi Directorate  Government of Karnataka    1458

Bank Of India     19791

Bank of Baroda  1412

CSC e-Governance Services India Limited          209771

Canara Bank      867

Commissioner Nagaland          25

DC Aalo  126

DC ITANAGAR CAPITAL COMPLEX          38

DC LOHIT          119

DC NAMSAI       154

DC PAPUMPARE          15

DC Siang 38

DENA BANK      33869

DIT Lakshadweep          1

Department of Information Technology Govt of Jharkhand          464

Dept of ITC Govt of Rajasthan     13565

Director General Health Services Health Deptt  Haryana  3

Director Health and Family Welfare  UT  165

Directorate of Public Health and Family Welfare  Govt of Andhra Pradesh          23

Directorate of Woman and Child Development  Government of Himachal Pradesh          33

FCR Govt of Haryana     1823

FCS Govt of Punjab       173

Govt of Goa      857

Govt of Gujarat 13894

Govt of Himachal Pradesh          583

Govt of Karnataka          5330

Govt of Kerala   11937

Govt of Maharashtra     241

Govt of Sikkim - Dept of Econo  32

Govt of UT of Chandigarh          95

Govt. of Mizoram          3220

Govt. of Uttarkhand      44

IDBI Bank ltd    31

Information Technology & Communication Department  3958

Madhya Pradesh State Electronics Development Corporation Ltd.          17309

NSDL e-Governance Infrastructure Limited          54214

National Cooperative Consumers Federation Of India Limited     2590

Odisha Computer Application Center     1701

Punjab National Bank     1400

Punjab and Sind Bank     1543

RDD Govt of Tripura     606

Registrar General India BEL2     167

Registrar General India ECIL     757

Registrar General India Others   7

Registrar General of India ITI     55

Rural Development Department Bihar-1 640

Rural Development Dept  Govt. of Bihar 4145

Secretery IT J&K     110

State Bank of India     3422

Tamil Nadu eGovernance Agency     15468

U P Electronics Corporation Limited     293

U.P. Development Systems Corporation Ltd     4139

UIDAI-Registrar 19

UT Govt. Of Dadra & Nagar Haveli     46

UT Of Daman and Diu   50

UT of Puducherry     1

UTI Infrastructure Technology & Services Limited     2395

Union Bank     5536

Women and Child Development Govt. of Jharkhand     39

Time taken: 47.462 seconds, Fetched: 60 row(s)

# 4. Find the number of states, districts in each state and sub-districts in each district

select State,count(*) from aadhar_ex group by State;

OK

Andaman and Nicobar Islands    7

Andhra Pradesh4540

Arunachal Pradesh        632

Assam   2972

Bihar    81776

Chandigarh        199

Chhattisgarh     4617

Dadra and Nagar Haveli 107

Daman and Diu 99

Delhi    7247

Goa       799

Gujarat 24146

Haryana           5138

Himachal Pradesh	1283

Jammu and Kashmir	1331

Jharkhand	7423

Karnataka	15755

Kerala	12378

Lakshadweep	5

Madhya Pradesh	37360

Maharashtra	19783

Manipur	562

Meghalaya	259

Mizoram	3172

Nagaland	392

Odisha	11972

Others	12

Puducherry	85

Punjab	5888

Rajasthan	28659

Sikkim	48

Tamil Nadu	21196

Telangana	3768

Tripura	726

Uttar Pradesh	69476

Uttarakhand	6521

West Bengal	60485

Time taken: 23.244 seconds, Fetched: 37 row(s)

hive>

```
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1565323849569_0002, Tracking URL = http://quickstart.cloudera:8088/proxy/application_1565323849569_0002/
Kill Command = /usr/lib/hadoop/bin/hadoop job  -kill job_1565323849569_0002
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2019-08-08 23:57:14,886 Stage-1 map = 0%,  reduce = 0%
2019-08-08 23:57:22,290 Stage-1 map = 100%,  reduce = 0%, Cumulative CPU 2.48 sec
2019-08-08 23:57:28,648 Stage-1 map = 100%,  reduce = 100%, Cumulative CPU 3.94 sec
MapReduce Total cumulative CPU time: 3 seconds 940 msec
Ended Job = job_1565323849569_0002
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1  Reduce: 1   Cumulative CPU: 3.94 sec   HDFS Read: 46492720 HDFS Write: 586 SUCCESS
Total MapReduce CPU Time Spent: 3 seconds 940 msec
OK
Andaman and Nicobar Islands     7
Andhra Pradesh  4540
Arunachal Pradesh       632
Assam   2972
Bihar   81776
Chandigarh      199
Chhattisgarh    4617
Dadra and Nagar Haveli  107
Daman and Diu   99
Delhi   7247
Goa     799
Gujarat 24146
Haryana 5138
Himachal Pradesh        1283
Jammu and Kashmir       1331
Jharkhand       7423
Karnataka       15755
Kerala  12378
Lakshadweep     5
```

hive> select State,count(distinct(district)) from aadhar_ex group by State;

OK

Andaman and Nicobar Islands    2

Andhra Pradesh13

Arunachal Pradesh          17

Assam    28

Bihar     38

Chandigarh         1

Chhattisgarh      30

Dadra and Nagar Haveli 1

Daman and Diu  2

Delhi      9

Goa      2

Gujarat 33

Haryana            21

Himachal Pradesh          11

Jammu and Kashmir        22

Jharkhand          24

| | |
|---|---|
| Karnataka | 30 |
| Kerala | 14 |
| Lakshadweep | 1 |
| Madhya Pradesh | 50 |
| Maharashtra | 36 |
| Manipur | 9 |
| Meghalaya | 8 |
| Mizoram | 8 |
| Nagaland | 11 |
| Odisha | 30 |
| Others | 1 |
| Puducherry | 2 |
| Punjab | 22 |
| Rajasthan | 33 |
| Sikkim | 4 |
| Tamil Nadu | 32 |
| Telangana | 10 |
| Tripura | 8 |
| Uttar Pradesh | 75 |
| Uttarakhand | 13 |
| West Bengal | 19 |

Time taken: 22.608 seconds, Fetched: 37 row(s)

```
Total MapReduce CPU Time Spent: 4 seconds 280 msec
OK
Andaman and Nicobar Islands      2
Andhra Pradesh  13
Arunachal Pradesh       17
Assam   28
Bihar   38
Chandigarh      1
Chhattisgarh    30
Dadra and Nagar Haveli  1
Daman and Diu   2
Delhi   9
Goa     2
Gujarat 33
Haryana 21
Himachal Pradesh        11
Jammu and Kashmir       22
Jharkhand       24
Karnataka       30
Kerala  14
Lakshadweep     1
Madhya Pradesh  50
Maharashtra     36
Manipur 9
Meghalaya       8
Mizoram 8
Nagaland        11
Odisha  30
Others  1
Puducherry      2
Punjab  22
Rajasthan       33
Sikkim  4
```

select district,count(distinct(sub_district)) from aadhar_ex group by district;

| | |
|---|---|
| Dhamtari | 4 |
| Dhanbad | 9 |
| Dhar | 7 |
| Dharmapuri | 5 |
| Dharwad | 5 |
| Dhemaji | 5 |
| Dhenkanal | 10 |
| Dholpur | 5 |
| Dhubri | 8 |
| Dhule | 4 |
| Dibrugarh | 8 |
| Dimapur | 6 |
| Dindigul | 9 |
| Dindori | 2 |
| Diu | 1 |
| Doda | 4 |

| | |
|---|---|
| Ganganagar | 9 |
| Ganjam | 22 |
| Garhwa | 14 |
| Gariyaband | 6 |
| Gautam Buddha Nagar | 4 |
| Gaya | 26 |
| Ghaziabad | 2 |
| Ghazipur | 4 |
| Gir Somnath | 6 |
| Giridih | 13 |
| Goalpara | 4 |
| Godda | 8 |
| Golaghat | 5 |
| Gomati | 6 |
| Gonda | 3 |
| Gondiya | 8 |
| Gopalganj | 14 |
| Gorakhpur | 6 |
| Gumla | 11 |
| Guna | 7 |
| Guntur | 53 |
| Gurdaspur | 3 |
| Gurgaon | 5 |
| Gwalior | 3 |
| Hailakandi | 4 |
| Hamirpur | 15 |
| Hanumangarh | 7 |
| Hapur | 3 |
| Harda | 3 |
| Hardoi | 5 |
| Haridwar | 3 |

Jamui    10

Janjgir-champa  10

Jashpur 8

Jaunpur 5

Jehanabad    8

Jhabua  5

Jhajjar  3

Jhalawar    7

Jhansi    4

Jharsuguda    10

Jhunjhunun    6

Jind    4

Jodhpur    7

Jorhat  5

Junagadh    11

K.v. Rangareddy    32

Kabeerdham    4

Kachchh    10

Kaimur (Bhabua)    11

Kaithal  2

Kalaburagi    9

Kalahandi    13

Kamrup 10

Kamrup Metro  4

Kancheepuram  18

Kandhamal    16

Kangra  23

Kanker  6

Kannauj    2

Kanniyakumari  4

Kannur  3

Kanpur Dehat    5

Kanpur Nagar    3



## 6. Find out the names of private agencies for each state.

select distinct(state),enrolment_Agency from aadhar_ex;

West Bengal     United Telecoms Ltd

West Bengal     United Telecoms e-Services Pvt Ltd

West Bengal     Urmila Info solution

West Bengal     Utility Forms Pvt Ltd

West Bengal     VAP INFOSOLUTIONS

West Bengal     VEETECHNOLOGIES PVT. LTD

West Bengal     VISION COMPTECH INTEGRATOR LTD

West Bengal     Vakrangee Softwares Limited

West Bengal     Vayam technologies Ltd

West Bengal     Vedavaag Systems Limited

West Bengal     Virinchi Technologies Ltd

West Bengal     WEBEL TECHNOLOGY LIMITED

West Bengal     Wipro Ltd

West Bengal        Zephyr System Pvt.Ltd.

Time taken: 23.996 seconds, Fetched: 2271 row(s)



# CHECKPOINT 3

## 8. Find top 3 states generating most number of Aadhaar cards?

hive> create table if not exists subaadhar as select state,sum(aadharGenerated) as aadharsum from aadhar_ex group by state;

select State,aadharsum from subaadhar order by aadharsum  desc limit 3;

OK

Bihar        162607

West Bengal        119901

Uttar Pradesh        103767

Time taken: 21.164 seconds, Fetched: 3 row(s)

## 9. Find top 3 private agencies generating the most number of Aadhar cards?

hive> create table if not exists subaadhar_pa as select enrolment_Agency,sum(aadhaargenerated) as aadharsum_pa from aadhar_ex group by enrolment_Agency;

select enrolment_Agency,aadharsum_pa from subaadhar_pa order by aadharsum_pa  desc limit 3;

OK

CSC SPV            173192

Wipro Ltd         39619

SREI INFRASTRUCTURE FINANCES L        26497

Time taken: 22.798 seconds, Fetched: 3 row(s)

## 10. Find the number of residents providing email, mobile number? (Hint: consider non-zero values.)

hive> select count(*) from aadhar_ex where residentsProvidingEmail<> 0 AND residentsProvidingMobileNumber<> 0;

OK

16951

Time taken: 23.218 seconds, Fetched: 1 row(s)

## 11. Find top 3 districts where enrolment numbers are maximum?

SELECT district,count(*) as cnt from aadhar_ex where enrolmentRejected = 0 group by district order by cnt desc limit 3;

OK

Barddhaman        6726

North 24 Parganas        6534

South 24 Parganas        5603

## 12. Find the no. of Aadhaar cards generated in each state?

select state,sum(aadhaarGenerated)from aadhar_ex group by state;

OK

Andaman and Nicobar Islands    5

Andhra Pradesh5798

Arunachal Pradesh         913

Assam   3213

Bihar     162607

Chandigarh         259

Chhattisgarh      6604

Dadra and Nagar Haveli 140

Daman and Diu 105

Delhi     8426

Goa       1167

Gujarat 34844

Haryana            6804

Himachal Pradesh    1547

Jammu and Kashmir    1234

Jharkhand    9868

Karnataka    19764

Kerala    15143

Lakshadweep    4

Madhya Pradesh    53276

Maharashtra    26085

Manipur    1323

Meghalaya    277

Mizoram    6279

Nagaland    545

Odisha    18182

Others    12

Puducherry    83

Punjab    6506

Rajasthan    39570

Sikkim    50

Tamil Nadu    32485

Telangana    5018

Tripura    908

Uttar Pradesh    103767

Uttarakhand    13227

West Bengal    119901

Time taken: 22.033 seconds, Fetched: 37 row(s)

hive>

# CHECKPOINT 4

## 13. Create a data frame using the file and provide its summary.

scala> aadharDF.describe()

res5: org.apache.spark.sql.DataFrame = [summary: string, age: string, aadharGenerated: string, enrolmentRejected: string, residentsProvidingEmail: string, residentsProvidingMobileNumber: string]



## 14. Write a command to see the correlation between "age" and "mobile_number"?

select corr(age,residentsProvidingMobileNumber)from aadhar_ex;

OK

-0.11754461896889339

Time taken: 24.314 seconds, Fetched: 1 row(s)

## 15. Find the number of unique pincodes in the data?

SELECT distinct(pin_code) from aadhar_ex;

854338

854339

854340

855101

855102

855105

855106

855107

855108

855113

855114

855115

855116

855117

855456

Others

Time taken: 24.98 seconds, Fetched: 17756 row(s)

## 16. Find the number of Aadhaar registrations rejected in Uttar Pradesh and Maharashtra?

select State,sum(enrolmentRejected) from aadhar_ex where state in ('Maharashtra','Uttar Pradesh') group by state;

OK

Maharashtra     1818

Uttar Pradesh   5286

Time taken: 22.562 seconds, Fetched: 2 row(s)

# CHECKPOINT 5

17. The top 3 states where the percentage of Aadhaar cards being generated for males is the highest.

hive> select state,(sum(aadhaarGenerated)*100)/(sum(aadhaarGenerated+enrolmentRejected))as male_members from aadhar_ex where gender='M' group by state order by male_members desc limit 3;

OK

Andaman and Nicobar Islands     100.0

Others     100.0

Lakshadweep     100.0

Time taken: 43.643 seconds, Fetched: 3 row(s)

18. In each of these 3 states, identify the top 3 districts where the percentage of Aadhaar cards being rejected for females is the highest.

hive> select district,(sum(enrolmentRejected)*100/(sum(aadhaarGenerated+enrolmentRejected))) as female_rejections from aadhar_ex where gender = 'F' and state in('Andaman and Nicobar Islands','Others','Lakshadweep')group by district order by female_rejections desc limit 3;

OK

Lakshadweep     100.0

South Andaman50.0

North And Middle Andaman      33.333333333333336

Time taken: 45.468 seconds, Fetched: 3 row(s)

## 19. The top 3 states where the percentage of Aadhaar cards being generated for females is the highest.

select state,(sum(aadhaarGenerated)*100)/(sum(aadhaarGenerated+enrolmentRejected))as female_num from aadhar_ex where gender='F' group by state order by female_num desc limit 3;
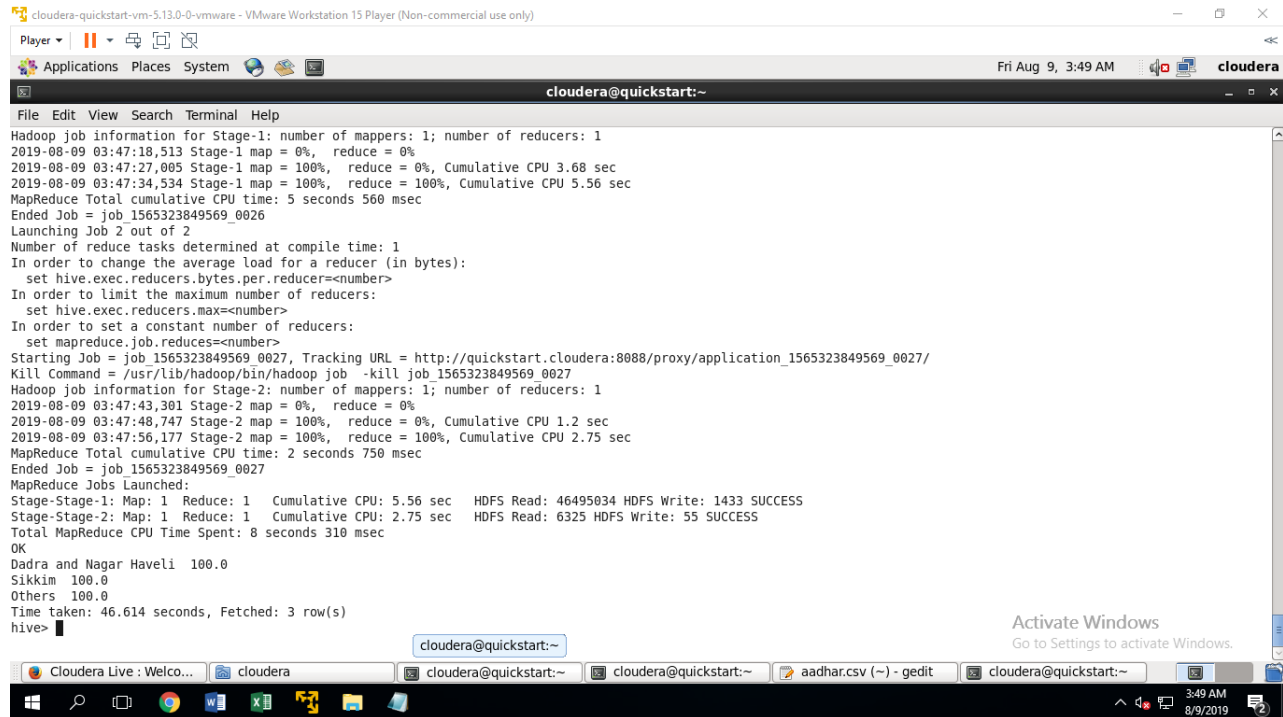
OK

Dadra and Nagar Haveli 100.0

Sikkim     100.0

Others     100.0

Time taken: 46.614 seconds, Fetched: 3 row(s)

```
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2019-08-09 03:47:18,513 Stage-1 map = 0%,  reduce = 0%
2019-08-09 03:47:27,005 Stage-1 map = 100%,  reduce = 0%, Cumulative CPU 3.68 sec
2019-08-09 03:47:34,534 Stage-1 map = 100%,  reduce = 100%, Cumulative CPU 5.56 sec
MapReduce Total cumulative CPU time: 5 seconds 560 msec
Ended Job = job_1565323849569_0026
Launching Job 2 out of 2
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1565323849569_0027, Tracking URL = http://quickstart.cloudera:8088/proxy/application_1565323849569_0027/
Kill Command = /usr/lib/hadoop/bin/hadoop job  -kill job_1565323849569_0027
Hadoop job information for Stage-2: number of mappers: 1; number of reducers: 1
2019-08-09 03:47:43,301 Stage-2 map = 0%,  reduce = 0%
2019-08-09 03:47:48,747 Stage-2 map = 100%,  reduce = 0%, Cumulative CPU 1.2 sec
2019-08-09 03:47:56,177 Stage-2 map = 100%,  reduce = 100%, Cumulative CPU 2.75 sec
MapReduce Total cumulative CPU time: 2 seconds 750 msec
Ended Job = job_1565323849569_0027
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1  Reduce: 1   Cumulative CPU: 5.56 sec   HDFS Read: 46495034 HDFS Write: 1433 SUCCESS
Stage-Stage-2: Map: 1  Reduce: 1   Cumulative CPU: 2.75 sec   HDFS Read: 6325 HDFS Write: 55 SUCCESS
Total MapReduce CPU Time Spent: 8 seconds 310 msec
OK
Dadra and Nagar Haveli  100.0
Sikkim  100.0
Others  100.0
Time taken: 46.614 seconds, Fetched: 3 row(s)
hive>
```

20. In each of these 3 states, identify the top 3 districts where the percentage of Aadhaar cards being rejected for males is the highest.

select district,(sum(enrolmentrejected)*100)/(sum(aadhaarGenerated+enrolmentRejected))as male_rejections from aadhar_ex where gender='M' and state in ('Dadra and Nagar Haveli','Sikkim','Others') group by district order by male_rejections desc limit 3;

OK
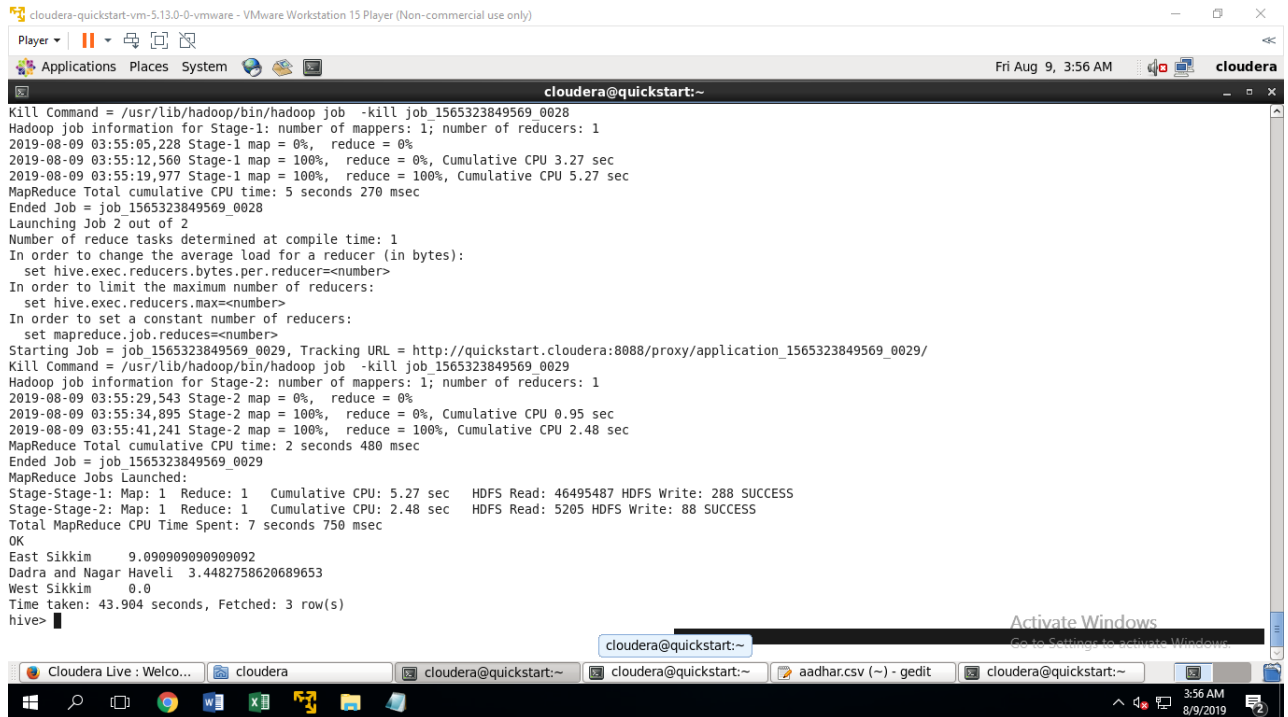East Sikkim          9.090909090909092
Dadra and Nagar Haveli   3.4482758620689653
West Sikkim          0.0
    Time taken: 43.904 seconds, Fetched: 3 row(s)

## 21. The summary of the acceptance percentage of all the Aadhaar cards applications by bucketing the age group into 10 buckets.

hive> create table if not exists aadhar_bucket(registrar String,enrolment_Agency String,State String,district String,sub_District String,pin_code String,gender String,age int,aadhaarGenerated int,enrolmentRejected int,residentsProvidingEmail int,residentsProvidingMobileNumber int)

> clustered by(age)into 10 buckets

> row format delimited fields terminated by ','

> stored as textfile

> TBLPROPERTIES('serialization.null.format'='',

> 'skip.header.line.count'='1');

OK

Time taken: 0.059 seconds

insert into aadhar_bucket select * from aadhar_ex;

select (sum(aadhaarGenerated)*100)/sum(aadhaarGenerated+enrolmentRejected)from aadhar_bucket;

OK

94.81863336350962

Time taken: 24.878 seconds, Fetched: 1 row(s)