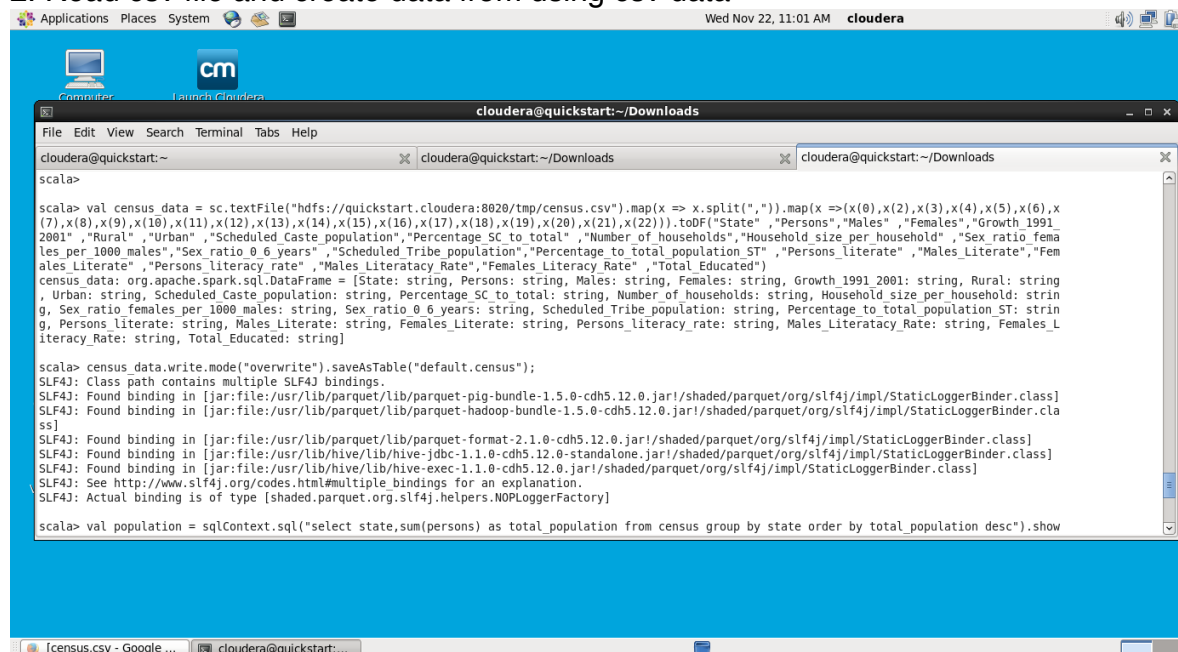Census data analysis
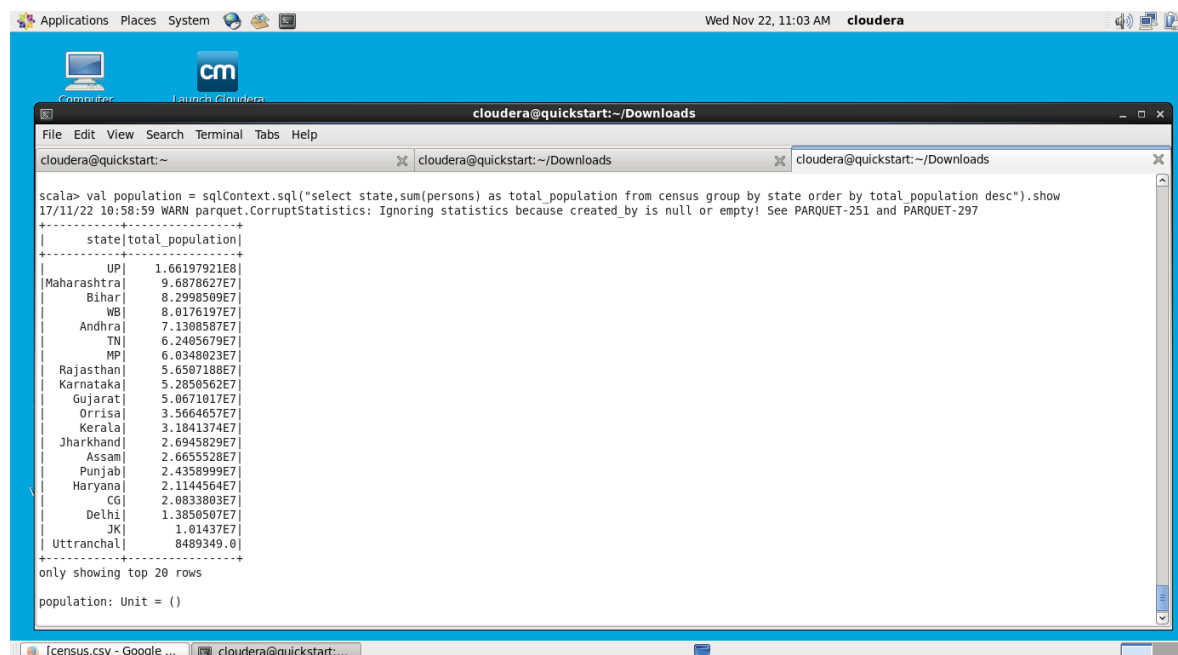
1. Copy data to HDFS directory using Hadoop dfs -put command
2. Read csv file and create data from using csv data



3. Find out the state wise population and order by state



Find out the Growth Rate of Each State Between 1991-2001

```
cloudera@quickstart:~/Downloads
File  Edit  View  Search  Terminal  Tabs  Help

cloudera@quickstart:~          cloudera@quickstart:~/Downloads          cloudera@quickstart:~/Downloads

scala> val growth_rate = sqlContext.sql("select state,avg(Growth_1991_2001) as total_growth from census group by state").show
+---------------+------------------+
|          state|      total_growth|
+---------------+------------------+
|    Maharashtra|19.607142857142865|
|             TN|10.127666666666668|
|         Gujarat|           20.8248|
|         Orrisa|15.551379310344826|
|         Sikkim|31.834999999999997|
|             AN|            18.665|
|      Chandigarh|             40.33|
|           Bihar|28.605945945945955|
|             HP| 17.53083333333333|
|             UP| 25.70228571428572|
|ArunachalPradesh| 25.46999999999999|
|         Tripura|15.405000000000001|
|            D_N_H|              59.2|
|      Uttranchal|17.092307692307692|
|         Haryana|27.816842105263152|
|              CG|17.506249999999998|
|              WB|18.424999999999997|
|         Manipur|29.240000000000002|
|              JK|28.785714285714285|
|      Lakshdweep|             17.19|
+---------------+------------------+
only showing top 20 rows

growth_rate: Unit = ()

scala>
```

# Find the literacy rate of each state

```
cloudera@quickstart:~/Downloads
File  Edit  View  Search  Terminal  Tabs  Help

cloudera@quickstart:~          cloudera@quickstart:~/Downloads          cloudera@quickstart:~/Downloads

scala> val literacy = sqlContext.sql("select state,avg(Persons_literacy_rate) from census group by state").show
+---------------+------------------+
|          state|                _c1|
+---------------+------------------+
|    Maharashtra| 74.55342857142857|
|             TN| 72.94266666666665|
|         Gujarat| 67.07480000000001|
|         Orrisa| 59.97965517241381|
|         Sikkim|           66.9975|
|             AN| 77.41999999999999|
|      Chandigarh|             81.94|
|           Bihar| 46.42135135135135|
|             HP| 75.50833333333333|
|             UP| 56.01057142857144|
|ArunachalPradesh| 53.166923076923084|
|         Tripura| 70.27000000000001|
|            D_N_H|             57.63|
|      Uttranchal| 72.01769230769231|
|         Haryana| 68.24473684210527|
|              CG| 63.02312499999999|
|              WB|             66.07|
|         Manipur|           68.6125|
|              JK| 54.867142857142845|
|      Lakshdweep|             86.66|
+---------------+------------------+
only showing top 20 rows

literacy: Unit = ()

scala>
```

# Find out the States with More Female Population

Access documents, folders and network places

Computer    cm Launch Cloudera

**cloudera@quickstart:~/Downloads**

File  Edit  View  Search  Terminal  Tabs  Help

cloudera@quickstart:~    cloudera@quickstart:~/Downloads    cloudera@quickstart:~/Downloads

```
scala> val female_pop = sqlContext.sql("select state, sum(Males)-sum(Females) from census group by state").show
+---------------+---------+
|          state|       _c1|
+---------------+---------+
|    Maharashtra|3922565.0|
|             TN| 396139.0|
|        Gujarat|2100137.0|
|         Orrisa| 482015.0|
|         Sikkim|  36117.0|
|             AN|  29792.0|
|      Chandigarh| 113241.0|
|           Bihar|3489081.0|
|              HP|  97980.0|
|              UP|8932817.0|
|ArunachalPradesh|  61914.0|
|          Tripura|  85247.0|
|            D_N_H|  22842.0|
|       Uttranchal| 162499.0|
|          Haryana|1583342.0|
|               CG| 114633.0|
|               WB|2755773.0|
|           Manipur|  20533.0|
|               JK| 578152.0|
|        Lakshdweep|   1612.0|
+---------------+---------+
only showing top 20 rows

female_pop: Unit = ()

scala>
```

[census.csv - Google ...    cloudera@quickstart:...

---

## Find out the Percentage of Population in Every State

Computer    cm Launch Cloudera

**cloudera@quickstart:~/Downloads**

File  Edit  View  Search  Terminal  Tabs  Help

cloudera@quickstart:~    cloudera@quickstart:~/Downloads    cloudera@quickstart:~/Downloads

```
scala>

scala> val percenet_pop = sqlContext.sql("select state, (sum(persons) * 100.0) / SUM(sum(persons)) over() as percent_pop_by_state from census group by state").show
17/11/22 11:10:56 WARN execution.Window: No Partition Defined for Window operation! Moving all data to a single partition, this can cause serious performance degradation.
+---------------+--------------------+
|          state|percent_pop_by_state|
+---------------+--------------------+
|    Maharashtra|    9.475494209385522|
|             TN|    6.103767861999858|
|        Gujarat|    4.956025317815201|
|         Orrisa|    3.488284891601744|
|         Sikkim|  0.05289949576432755|
|             AN|  0.03483447606726582|
|      Chandigarh|  0.08808921009243792|
|           Bihar|    8.117909138174843|
|              HP|    0.5944665819347776|
|              UP|    16.25546817511578|
|ArunachalPradesh|  0.10738993468694186|
|          Tripura|  0.31290729895613395|
|            D_N_H|  0.02156566193106157|
|       Uttranchal|    0.8303253233652121|
|          Haryana|    2.0681052152192616|
|               CG|    2.0377103371415317|
|               WB|    7.841864753141607|
|           Manipur|    0.19662075848548596|
|               JK|    0.9921339059826262|
|        Lakshdweep|0.005932048601382...|
+---------------+--------------------+
only showing top 20 rows
```

[census.csv - Google ...    cloudera@quickstart:...