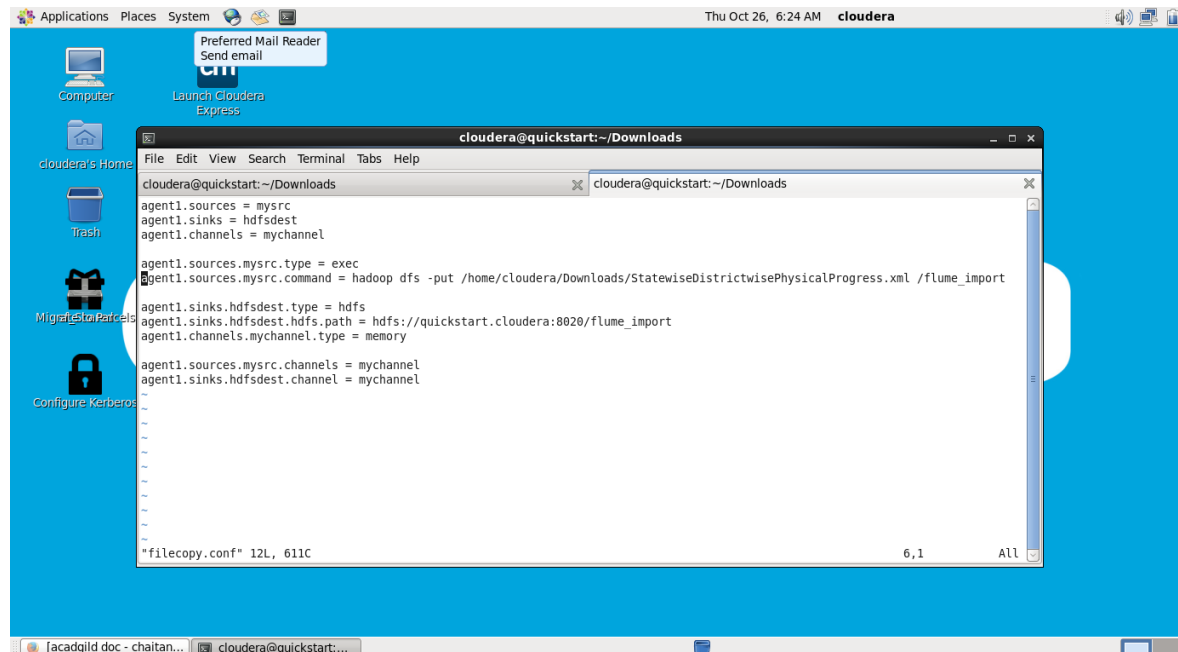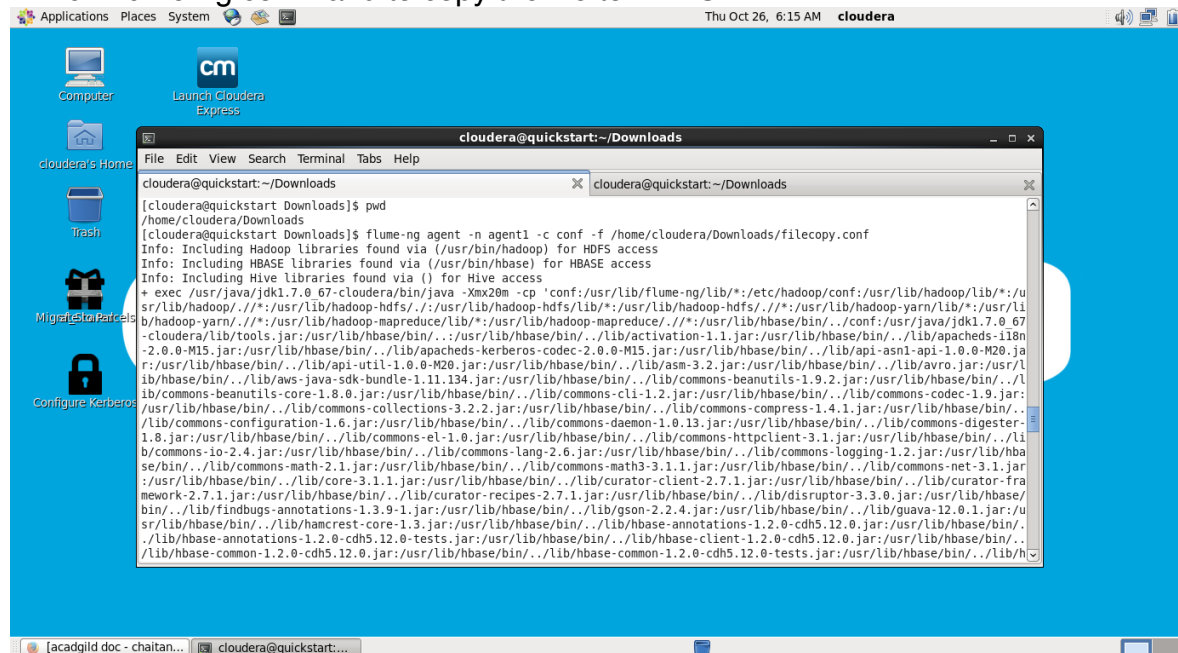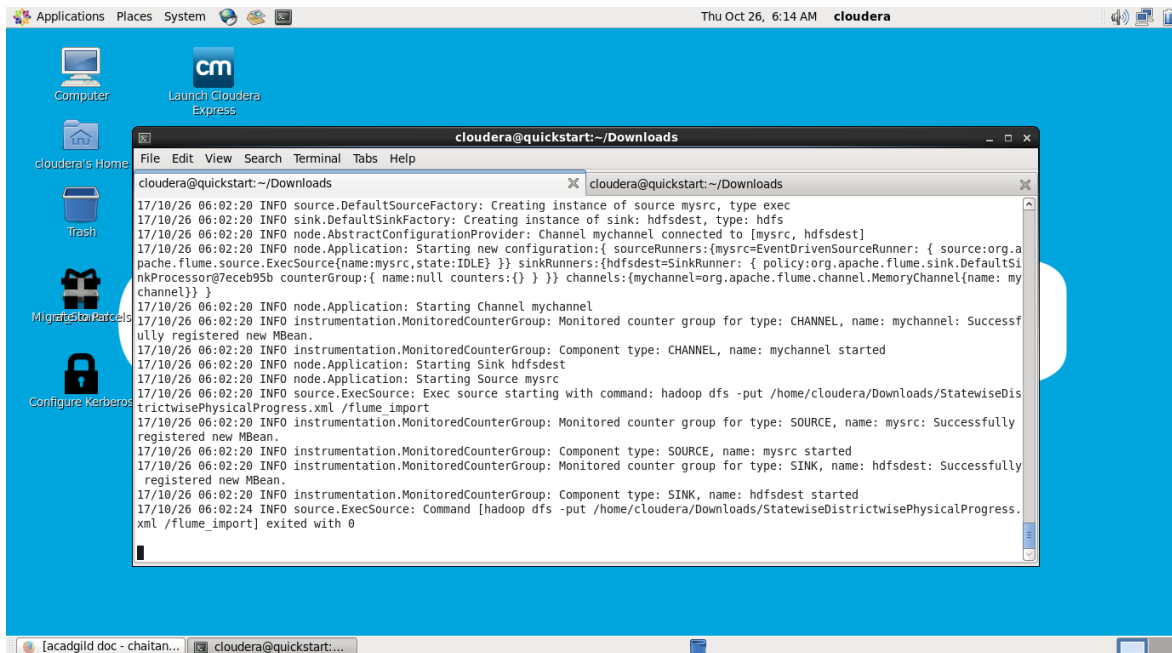Task 1:
The FLUME job which will format the data and place the data to HDFS

1. Conf file to download dataset from local file system to HDFS flume:
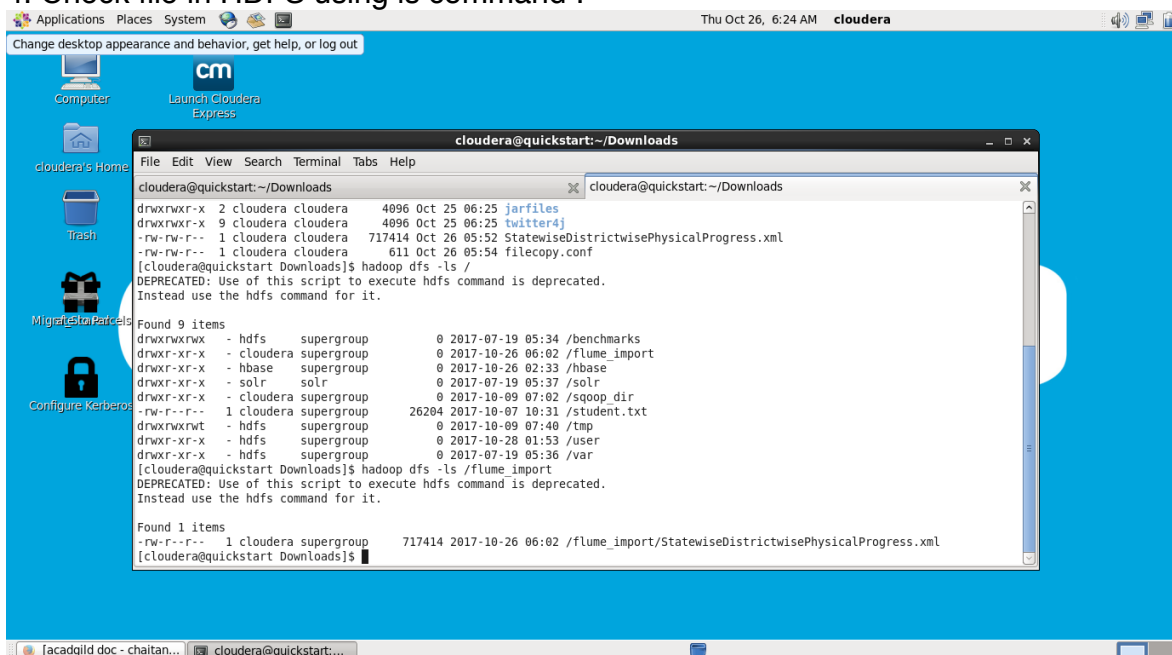


2. Run flume-ng command to copy the file to HDFS.

## 4. Check file in HDFS using ls command :



Task 2:
Pig/MapReduce job for parsing the XML data.

Pig Script:

```
REGISTER '/home/cloudera/Downloads/jarfiles/piggybank-0.17.0.jar'
DEFINE XPath org.apache.pig.piggybank.evaluation.xml.XPath();
A = LOAD '/flume_import/StatewiseDistrictwisePhysicalProgress.xml' using
org.apache.pig.piggybank.storage.XMLLoader('row') as (x:chararray);
B = FOREACH A GENERATE XPath(x,'row/Project_Objectives_IHHL_BPL'),
XPath(x,'row/Project_Performance-IHHL_BPL');
dump B;
```

Execution:
Pig <pig_file_name>

Output:



Task 3:
Create Pig scripts/MapReduce jobs to analyze the data

Find out the districts who achieved 100 percent objective in BPL cards
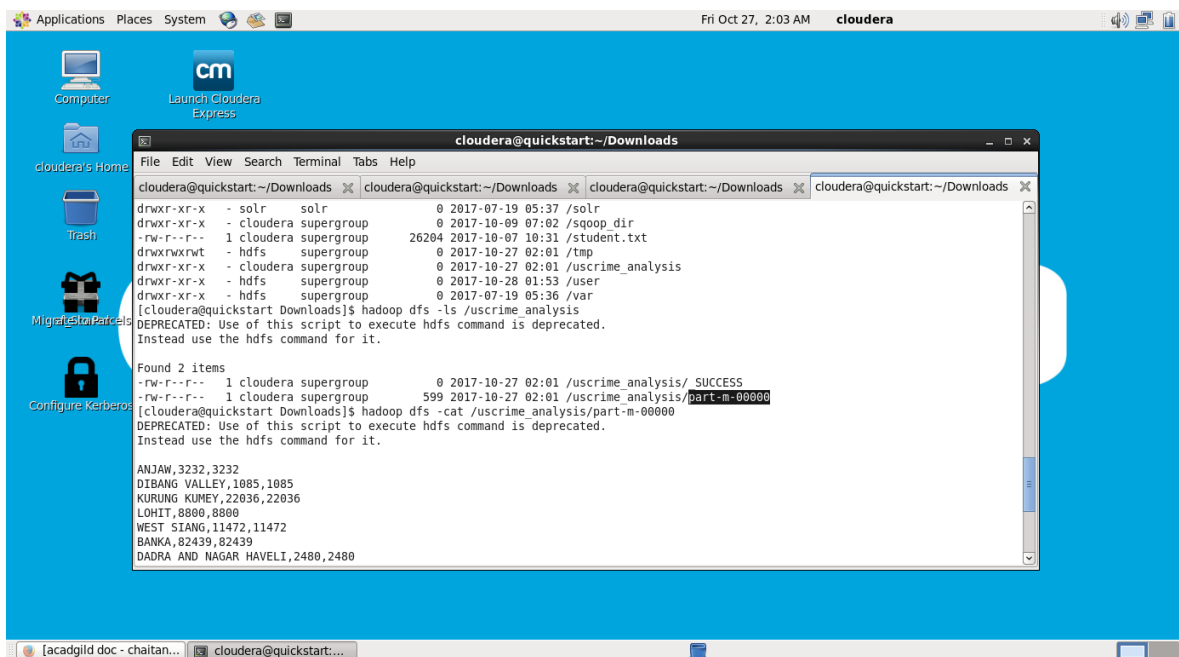
Pig script:
REGISTER '/home/cloudera/Downloads/jarfiles/piggybank-0.17.0.jar'
DEFINE XPath org.apache.pig.piggybank.evaluation.xml.XPath();
A = LOAD '/flume_import/StatewiseDistrictwisePhysicalProgress.xml' using
org.apache.pig.piggybank.storage.XMLLoader('row') as (x:chararray);
B = FOREACH A GENERATE XPath(x,'row/District_Name') as district ,XPath(x,'row/
Project_Objectives_IHHL_BPL') as BPL_Objective, XPath(x,'row/
Project_Objectives_IHHL_TOTAL') as BPL_Objective_total;
C = filter B by (((int)BPL_Objective * 100)/(int)BPL_Objective_total) == 100;
STORE C INTO 'hdfs://quickstart.cloudera:8020/uscrime_analysis' USING
PigStorage (',');
dump C;

Execution:
Pig <file_name_path>

Output:

Export the results to mysql using sqoop:

1. Create table in mysql

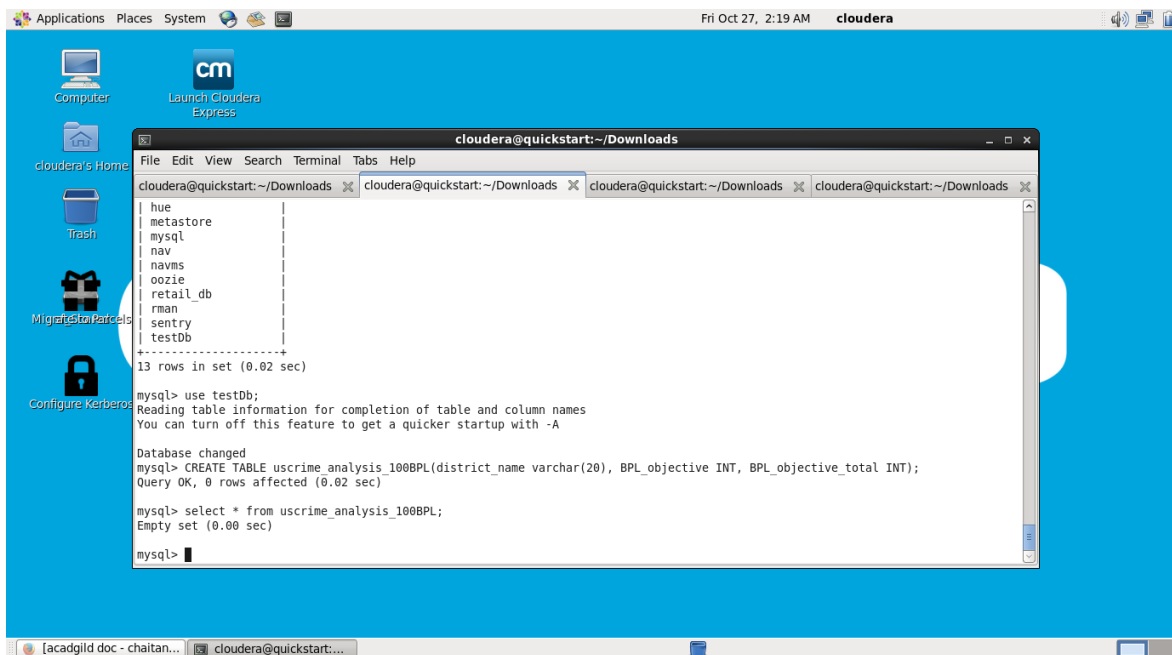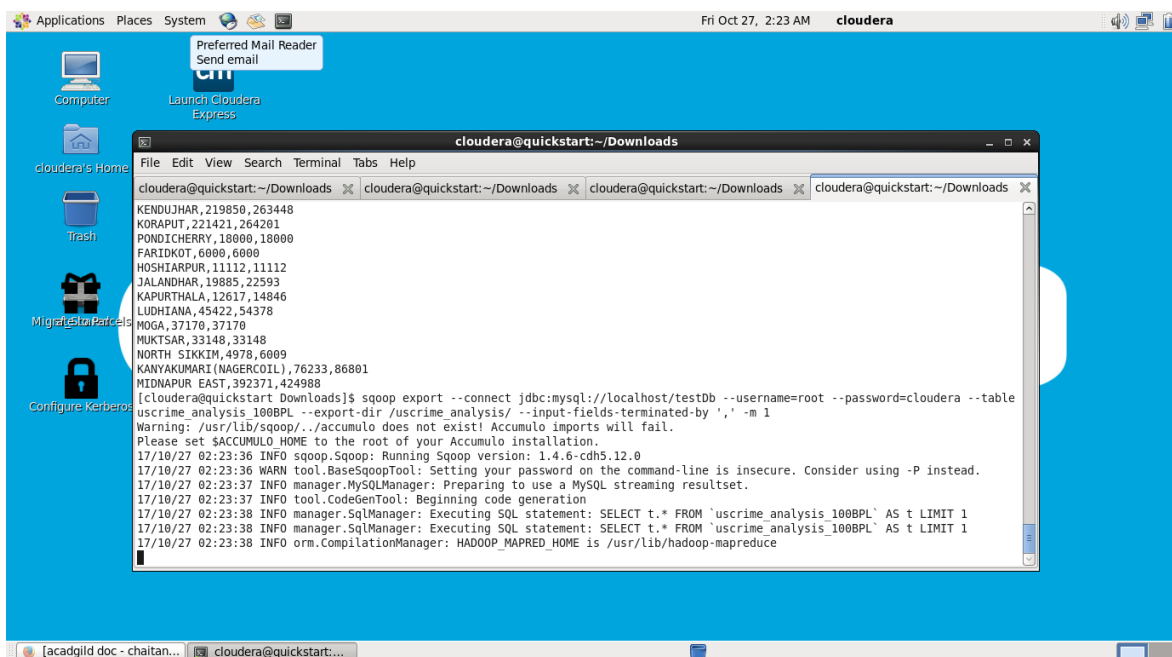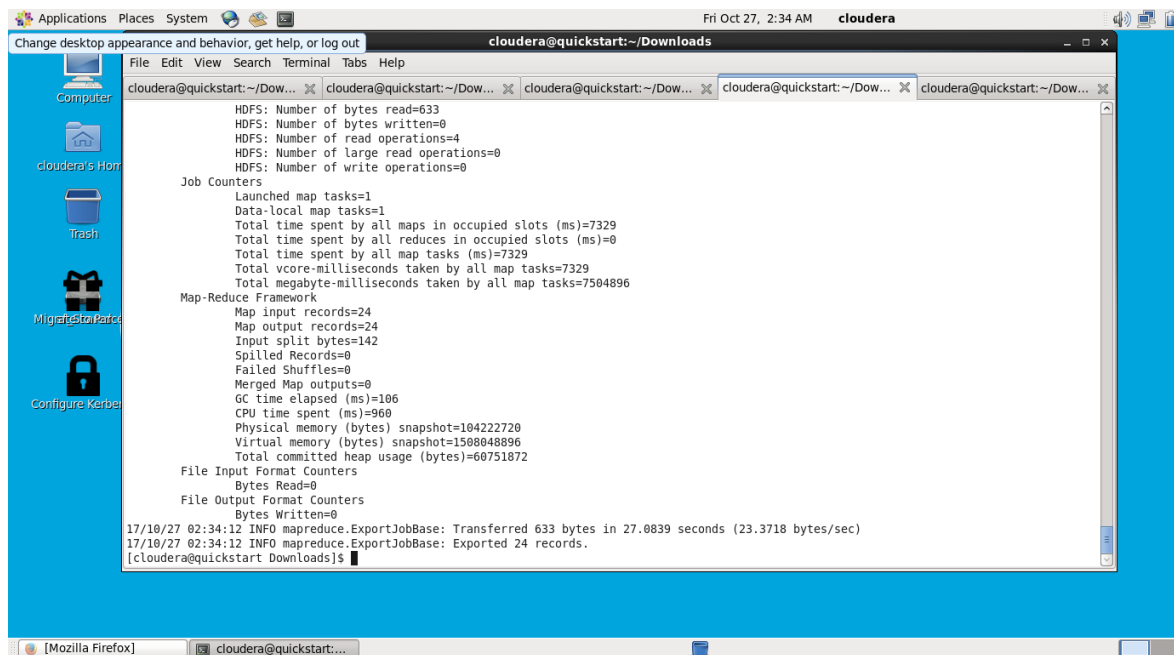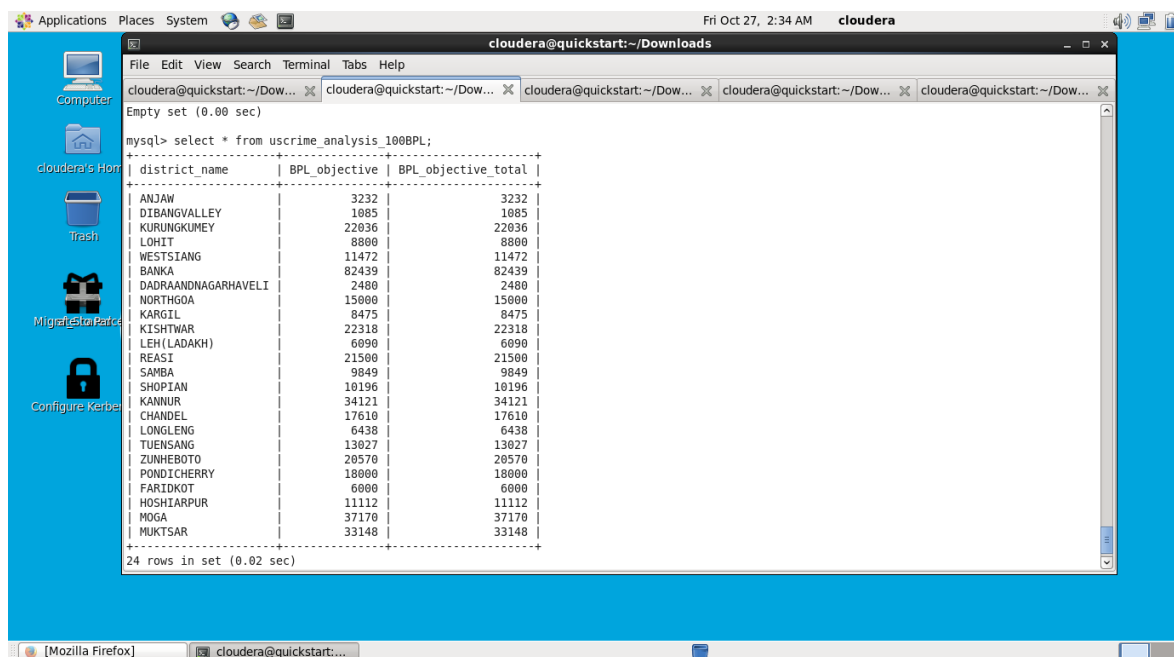## 2. Scoop command to export the data from HDFS to mysql



## 3. once the job completed, check for the success command:

## 4. Check the exported data in mysql select statement:

```
Empty set (0.00 sec)

mysql> select * from uscrime_analysis_100BPL;
+---------------------+---------------+---------------------+
| district_name       | BPL_objective | BPL_objective_total |
+---------------------+---------------+---------------------+
| ANJAW               |          3232 |                3232 |
| DIBANGVALLEY        |          1085 |                1085 |
| KURUNGKUMEY         |         22036 |               22036 |
| LOHIT               |          8800 |                8800 |
| WESTSIANG           |         11472 |               11472 |
| BANKA               |         82439 |               82439 |
| DADRAANDNAGARHAVELI |          2480 |                2480 |
| NORTHGOA            |         15000 |               15000 |
| KARGIL              |          8475 |                8475 |
| KISHTWAR            |         22318 |               22318 |
| LEH(LADAKH)         |          6090 |                6090 |
| REASI               |         21500 |               21500 |
| SAMBA               |          9849 |                9849 |
| SHOPIAN             |         10196 |               10196 |
| KANNUR              |         34121 |               34121 |
| CHANDEL             |         17610 |               17610 |
| LONGLENG            |          6438 |                6438 |
| TUENSANG            |         13027 |               13027 |
| ZUNHEBOTO           |         20570 |               20570 |
| PONDICHERRY         |         18000 |               18000 |
| FARIDKOT            |          6000 |                6000 |
| HOSHIARPUR          |         11112 |               11112 |
| MOGA                |         37170 |               37170 |
| MUKTSAR             |         33148 |               33148 |
+---------------------+---------------+---------------------+
24 rows in set (0.02 sec)
```

2.Write a Pig UDF to filter the districts which have reached 80% of objectives of BPL cards.

Pig script:
REGISTER '/home/cloudera/Downloads/jarfiles/piggybank-0.17.0.jar'
DEFINE XPath org.apache.pig.piggybank.evaluation.xml.XPath();
A = LOAD '/flume_import/StatewiseDistrictwisePhysicalProgress.xml' using
org.apache.pig.piggybank.storage.XMLLoader('row') as (x:chararray);
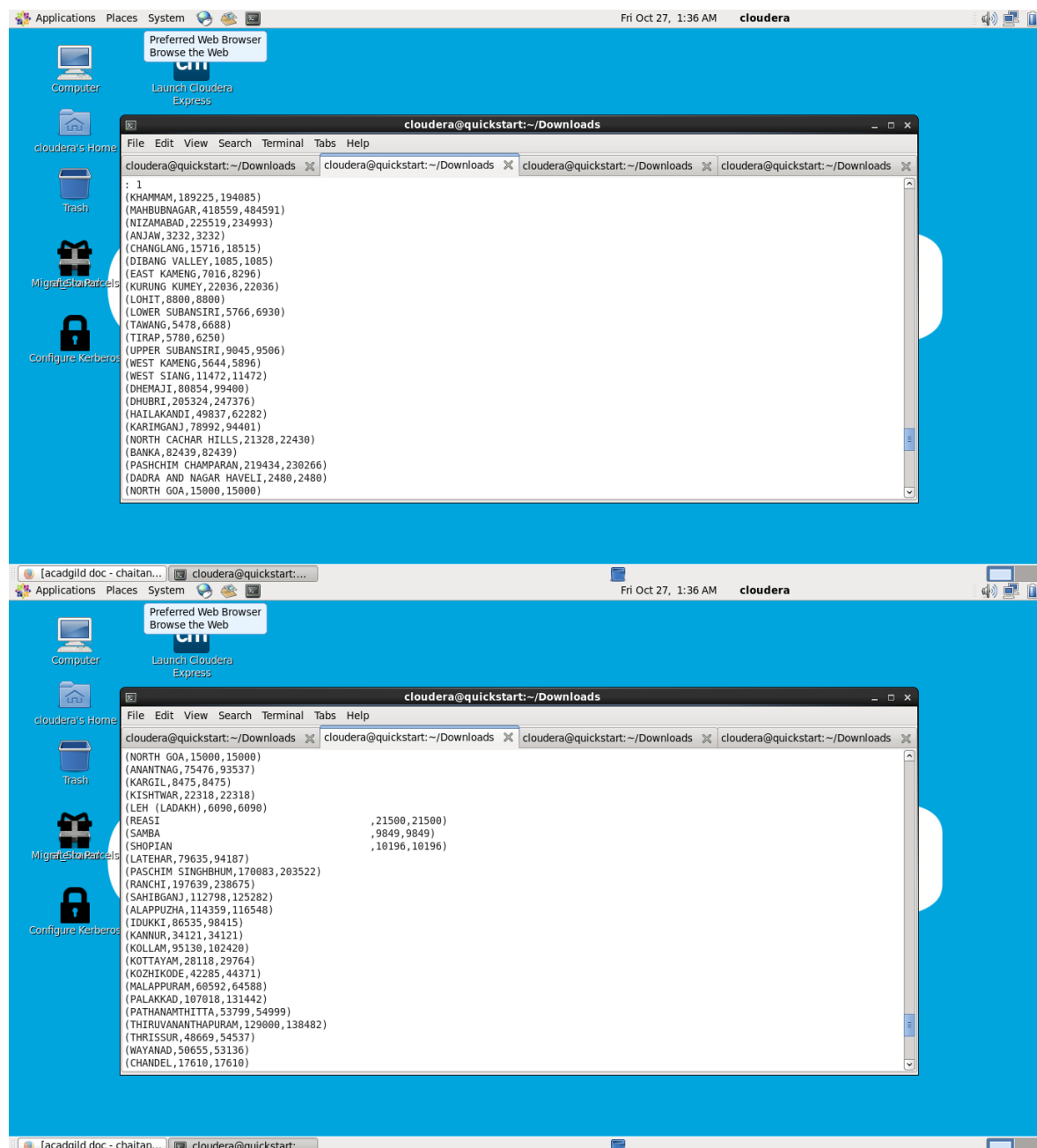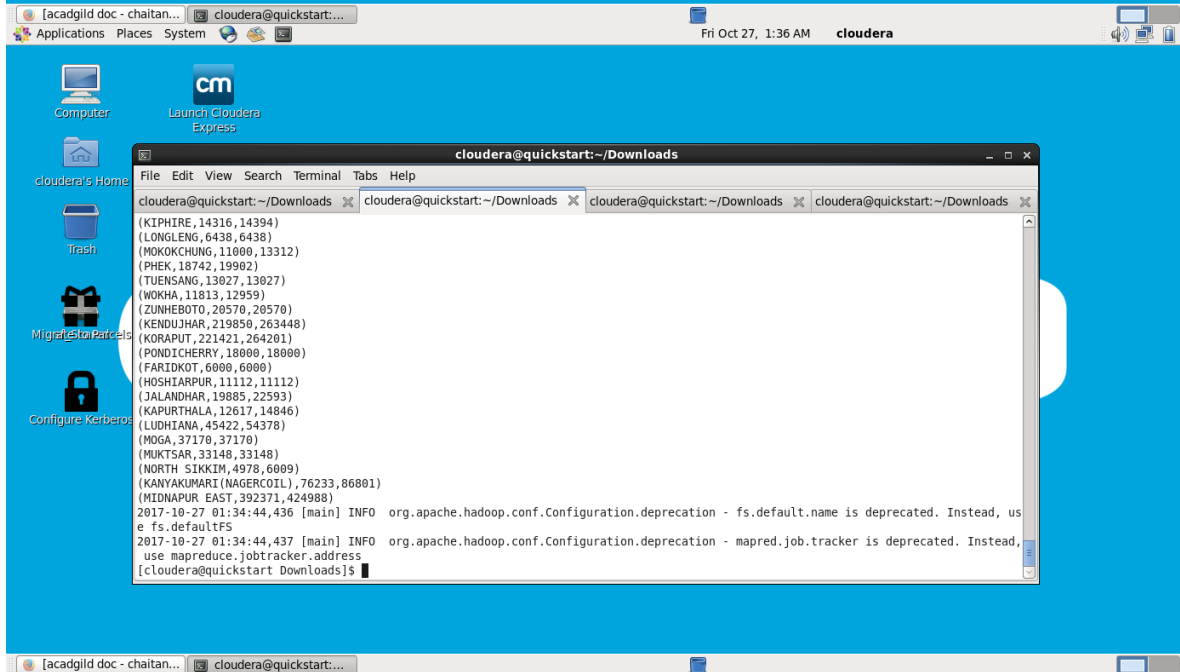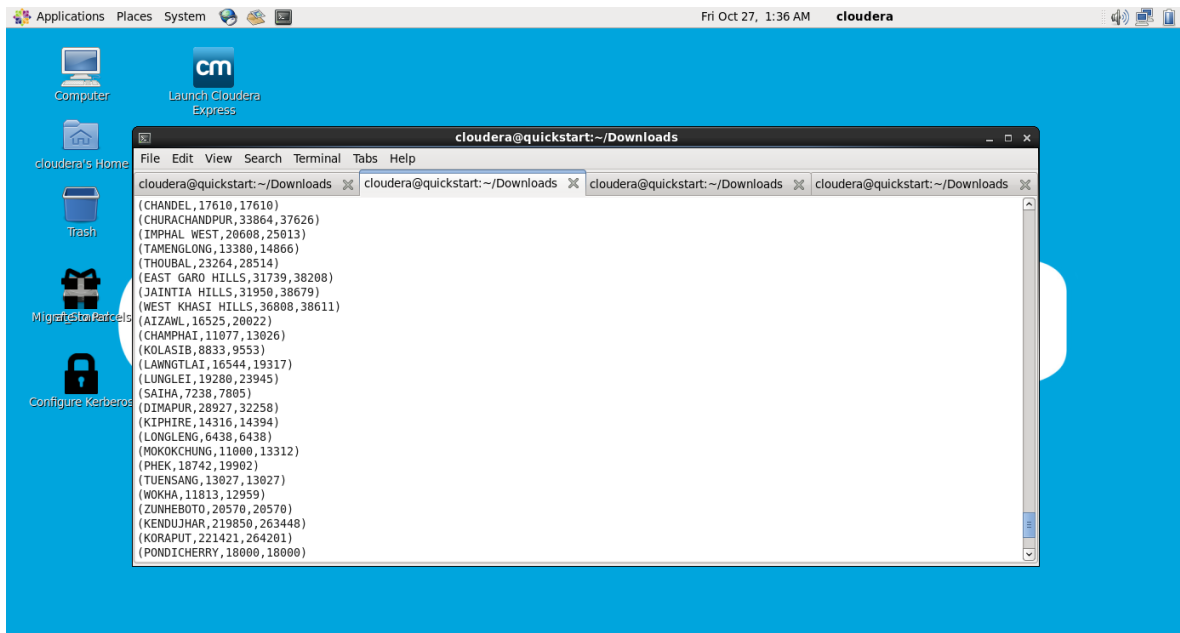B = FOREACH A GENERATE XPath(x,'row/District_Name') as district ,XPath(x,'row/

Project_Objectives_IHHL_BPL') as BPL_Objective, XPath(x,'row/
Project_Objectives_IHHL_TOTAL') as BPL_Objective_total;
C = filter B by (((int)BPL_Objective * 100)/(int)BPL_Objective_total) >= 80;
STORE C INTO 'hdfs://quickstart.cloudera:8020/uscrime_analysis_2' USING
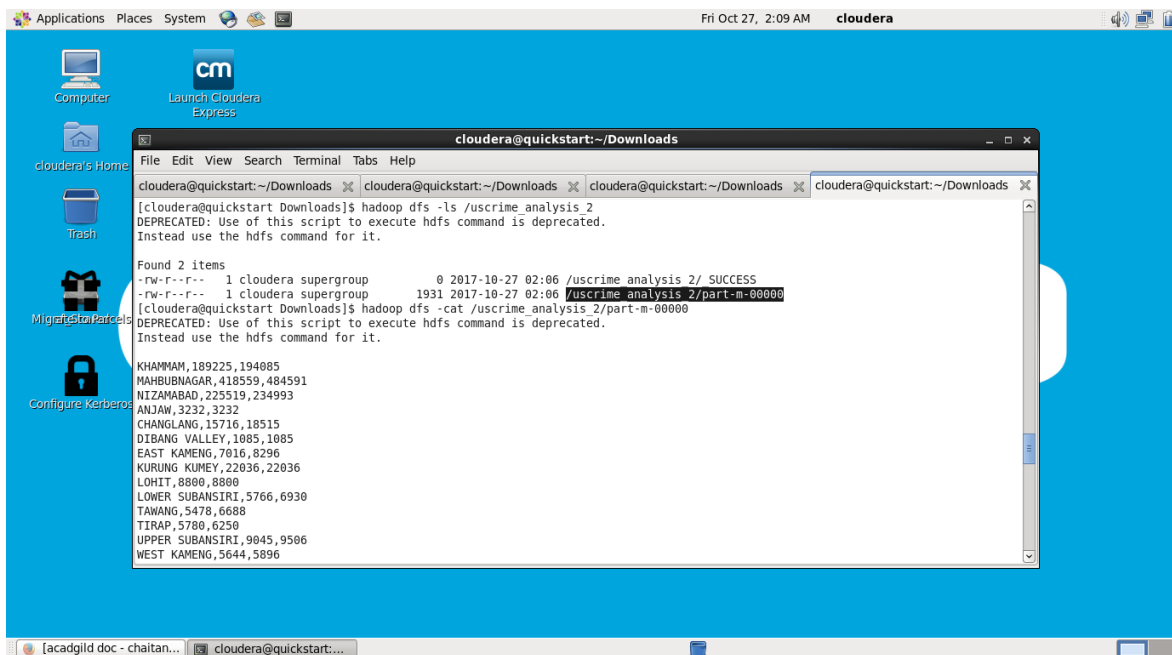PigStorage (',');
dump C;

Execution:
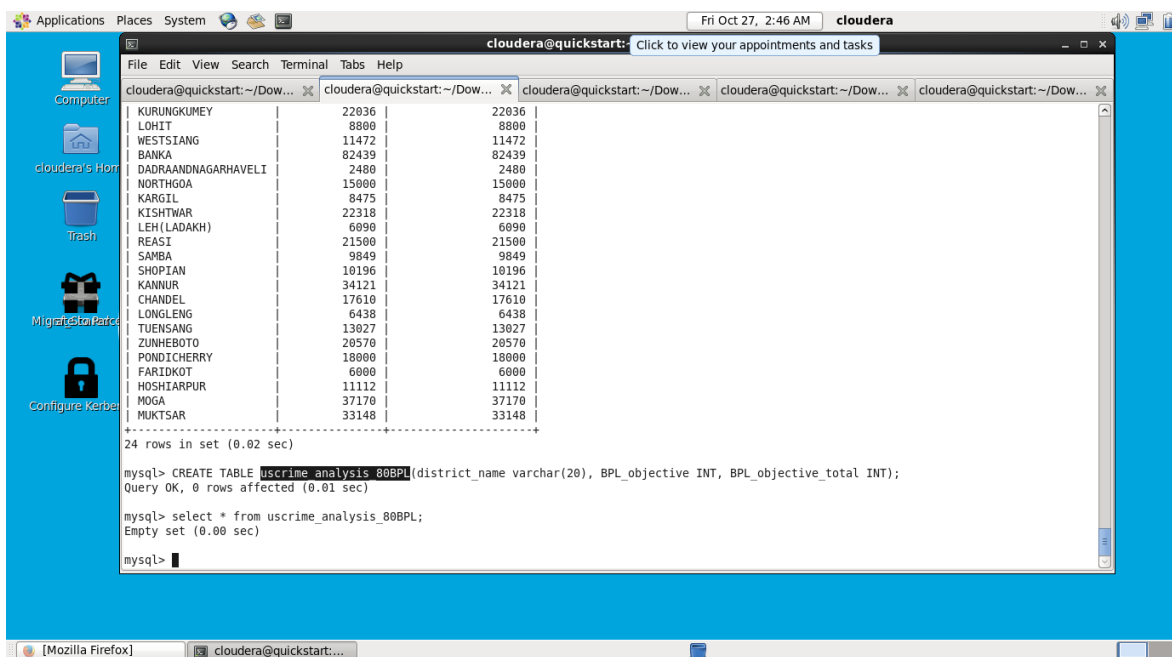Pig <pig_script_filename>

Output:

**cloudera@quickstart:~/Downloads**

File  Edit  View  Search  Terminal  Tabs  Help

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

```
(CHANDEL,17610,17610)
(CHURACHANDPUR,33864,37626)
(IMPHAL WEST,20608,25013)
(TAMENGLONG,13380,14866)
(THOUBAL,23264,28514)
(EAST GARO HILLS,31739,38208)
(JAINTIA HILLS,31950,38679)
(WEST KHASI HILLS,36808,38611)
(AIZAWL,16525,20022)
(CHAMPHAI,11077,13026)
(KOLASIB,8833,9553)
(LAWNGTLAI,16544,19317)
(LUNGLEI,19280,23945)
(SAIHA,7238,7805)
(DIMAPUR,28927,32258)
(KIPHIRE,14316,14394)
(LONGLENG,6438,6438)
(MOKOKCHUNG,11000,13312)
(PHEK,18742,19902)
(TUENSANG,13027,13027)
(WOKHA,11813,12959)
(ZUNHEBOTO,20570,20570)
(KENDUJHAR,219850,263448)
(KORAPUT,221421,264201)
(PONDICHERRY,18000,18000)
```

**cloudera@quickstart:~/Downloads**

File  Edit  View  Search  Terminal  Tabs  Help

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

```
(KIPHIRE,14316,14394)
(LONGLENG,6438,6438)
(MOKOKCHUNG,11000,13312)
(PHEK,18742,19902)
(TUENSANG,13027,13027)
(WOKHA,11813,12959)
(ZUNHEBOTO,20570,20570)
(KENDUJHAR,219850,263448)
(KORAPUT,221421,264201)
(PONDICHERRY,18000,18000)
(FARIDKOT,6000,6000)
(HOSHIARPUR,11112,11112)
(JALANDHAR,19885,22593)
(KAPURTHALA,12617,14846)
(LUDHIANA,45422,54378)
(MOGA,37170,37170)
(MUKTSAR,33148,33148)
(NORTH SIKKIM,4978,6009)
(KANYAKUMARI(NAGERCOIL),76233,86801)
(MIDNAPUR EAST,392371,424988)
2017-10-27 01:34:44,436 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, us
e fs.defaultFS
2017-10-27 01:34:44,437 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead,
 use mapreduce.jobtracker.address
[cloudera@quickstart Downloads]$
```

Export the results to mysql using sqoop:

1. Create table in mysql to store the data:



3. Using sqoop export data from HDFS to mysql using command:

## 3. Check the success message for the job completion:



## 4. Check the data in mysql table using command; 82 rows copied successfully.

cloudera@quickstart:~/Downloads

| cloudera@quickstart:~/Dow... | cloudera@quickstart:~/Dow... | cloudera@quickstart:~/Dow... | cloudera@quickstart:~/Dow... | cloudera@quickstart:~/Dow... |

```
| WESTKHASIHILLS          |    36808 |             38611 |
| AIZAWL                  |    16525 |             20022 |
| CHAMPHAI                |    11077 |             13026 |
| KOLASIB                 |     8833 |              9553 |
| LAWNGTLAI               |    16544 |             19317 |
| LUNGLEI                 |    19280 |             23945 |
| SAIHA                   |     7238 |              7805 |
| DIMAPUR                 |    28927 |             32258 |
| KIPHIRE                 |    14316 |             14394 |
| LONGLENG                |     6438 |              6438 |
| MOKOKCHUNG              |    11000 |             13312 |
| PHEK                    |    18742 |             19902 |
| TUENSANG                |    13027 |             13027 |
| WOKHA                   |    11813 |             12959 |
| ZUNHEBOTO               |    20570 |             20570 |
| KENDUJHAR               |   219850 |            263448 |
| KORAPUT                 |   221421 |            264201 |
| PONDICHERRY             |    18000 |             18000 |
| FARIDKOT                |     6000 |              6000 |
| HOSHIARPUR              |    11112 |             11112 |
| JALANDHAR               |    19885 |             22593 |
| KAPURTHALA              |    12617 |             14846 |
| LUDHIANA                |    45422 |             54378 |
| MOGA                    |    37170 |             37170 |
| MUKTSAR                 |    33148 |             33148 |
| NORTHSIKKIM             |     4978 |              6009 |
| KANYAKUMARI(NAGERCOIL)  |    76233 |             86801 |
| MIDNAPUREAST            |   392371 |            424988 |
+-------------------------+----------+-------------------+
82 rows in set (0.03 sec)

mysql>
```