

## Project 2:

### State-Wise Development Analysis In India

In this project we are performing the data parsing and data analysis using Pig and export the results into MYSQL using Sqoop.

#### Step 1:

1. Place the flume config file 'filecopy.conf' at location /home/acadgild/apache-flume-1.8.0-bin/conf

filecopy.conf

```
[acadgild@localhost ~]$ cat apache-flume-1.8.0-bin/conf/filecopy.conf
agent1.sources = mysrc
agent1.sinks = hdfsdest
agent1.channels = mychannel

agent1.sources.mysrc.type = exec
agent1.sources.mysrc.command = hadoop dfs -put /home/acadgild/StatewiseDistrictwisePhysicalProgress.xml /flume_import

agent1.sinks.hdfsdest.type = hdfs
agent1.sinks.hdfsdest.hdfs.path = hdfs://localhost:9000/flume_import

agent1.channels.mychannel.type = memory

agent1.sources.mysrc.channels = mychannel
agent1.sinks.hdfsdest.channel = mychannel
[acadgild@localhost ~]$
```

Copy dataset from local file system to HDFS using flume.

#### Command:

flume-ng agent -n agent1 -c conf -f /home/acadgild/apache-flume-1.8.0-bin/conf/filecopy.conf

```
[acadgild@localhost ~]$ flume-ng agent -n agent1 -c conf -f /home/acadgild/apache-flume-1.8.0-bin/conf/filecopy.conf
Info: Including Hadoop libraries found via (/usr/local/hadoop-2.6.0/bin/hadoop) for HDFS access
Info: Including HBASE libraries found via (/usr/local/hbase/bin/hbase) for HBASE access
Info: Including Hive libraries found via (/usr/local/hive) for Hive access
+ exec /usr/local/java/bin/java -Xmx20m -cp 'conf:/home/acadgild/apache-flume-1.8.0-bin/lib/*:/usr/local/hadoop-2.6.0/contrib/capacity-sc
heduler/*:/usr/local/hadoop-2.6.0/etc/hadoop:/usr/local/hadoop-2.6.0/share/hadoop/common/lib/*:/usr/local/hadoop-2.6.0/share/hadoop/c
ommon/*:/usr/local/hadoop-2.6.0/share/hadoop/hdfs:/usr/local/hadoop-2.6.0/share/hadoop/hdfs/lib/*:/usr/local/hadoop-2.6.0/share/hadoop/hd
fs/*:/usr/local/hadoop-2.6.0/share/hadoop/yarn/lib/*:/usr/local/hadoop-2.6.0/share/hadoop/yarn/*:/usr/local/hadoop-2.6.0/share/hadoop/map
reduce/lib/*:/usr/local/hadoop-2.6.0/share/hadoop/mapreduce/*:/usr/local/hbase/conf:/usr/local/java/lib/tools.jar:/usr/local/hbase:/usr/l
ocal/hbase/lib/activation-1.1.jar:/usr/local/hbase/lib/aopalliance-1.0.jar:/usr/local/hbase/lib/asm-3.1.jar:/usr/local/hbase/lib/avro-1.7
.4.jar:/usr/local/hbase/lib/commons-beanutils-1.7.0.jar:/usr/local/hbase/lib/commons-beanutils-core-1.8.0.jar:/usr/local/hbase/lib/common
s-cli-1.2.jar:/usr/local/hbase/lib/commons-codec-1.7.jar:/usr/local/hbase/lib/commons-collections-3.2.1.jar:/usr/local/hbase/lib/commons-
compress-1.4.1.jar:/usr/local/hbase/lib/commons-configuration-1.6.jar:/usr/local/hbase/lib/commons-daemon-1.0.13.jar:/usr/local/hbase/lib/co
mmons-digester-1.8.jar:/usr/local/hbase/lib/commons-el-1.0.jar:/usr/local/hbase/lib/commons-httpclient-3.1.jar:/usr/local/hbase/lib/co
mmons-io-2.4.jar:/usr/local/hbase/lib/commons-lang-2.6.jar:/usr/local/hbase/lib/commons-logging-1.1.1.jar:/usr/local/hbase/lib/commons-ma
th-2.1.jar:/usr/local/hbase/lib/commons-net-3.1.jar:/usr/local/hbase/lib/findbugs-annotations-1.3.9-1.jar:/usr/local/hbase/lib/gmbal-api-
only-3.0.0-b023.jar:/usr/local/hbase/lib/grizzlly-framework-2.1.2.jar:/usr/local/hbase/lib/grizzlly-http-2.1.2.jar:/usr/local/hbase/lib/gri
zzlly-http-server-2.1.2.jar:/usr/local/hbase/lib/grizzlly-http-servlet-2.1.2.jar:/usr/local/hbase/lib/grizzlly-rcm-2.1.2.jar:/usr/local/hbas
e/lib/guava-12.0.1.jar:/usr/local/hbase/lib/guice-3.0.jar:/usr/local/hbase/lib/guice-servlet-3.0.jar:/usr/local/hbase/lib/hadoop-annotati
ons-2.2.0.jar:/usr/local/hbase/lib/hadoop-auth-2.2.0.jar:/usr/local/hbase/lib/hadoop-client-2.2.0.jar:/usr/local/hbase/lib/hadoop-common-
2.2.0.jar:/usr/local/hbase/lib/hadoop-hdfs-2.2.0.jar:/usr/local/hbase/lib/hadoop-mapreduce-client-app-2.2.0.jar:/usr/local/hbase/lib/hadoo
p-mapreduce-client-common-2.2.0.jar:/usr/local/hbase/lib/hadoop-mapreduce-client-core-2.2.0.jar:/usr/local/hbase/lib/hadoop-mapreduce-cl
ient-jobclient-2.2.0.jar:/usr/local/hbase/lib/hadoop-mapreduce-client-shuffle-2.2.0.jar:/usr/local/hbase/lib/hadoop-yarn-api-2.2.0.jar:/u
sr/local/hbase/lib/hadoop-yarn-client-2.2.0.jar:/usr/local/hbase/lib/hadoop-yarn-common-2.2.0.jar:/usr/local/hbase/lib/hadoop-yarn-server
-common-2.2.0.jar:/usr/local/hbase/lib/hadoop-yarn-server-nodemanager-2.2.0.jar:/usr/local/hbase/lib/hamcrest-core-1.3.jar:/usr/local/hba
se/lib/hbase-annotations-0.98.14-hadoop2.jar:/usr/local/hbase/lib/hbase-checkstyle-0.98.14-hadoop2.jar:/usr/local/hbase/lib/hbase-client-
0.98.14-hadoop2.jar:/usr/local/hbase/lib/hbase-common-0.98.14-hadoop2.jar:/usr/local/hbase/lib/hbase-common-0.98.14-hadoop2-tests.jar:/us
r/local/hbase/lib/hbase-examples-0.98.14-hadoop2.jar:/usr/local/hbase/lib/hbase-hadoop2-compat-0.98.14-hadoop2.jar:/usr/local/hbase/lib/h
base-hadoop-compat-0.98.14-hadoop2.jar:/usr/local/hbase/lib/hbase-it-0.98.14-hadoop2.jar:/usr/local/hbase/lib/hbase-protocol-0.98.14-hadoop2.jar:/usr/local/hba
se/lib/hbase-resource-bundle-0.98.14-hadoop2.jar:/usr/local/hbase/lib/hbase-rest-0.98.14-hadoop2.jar:/usr/local/hbase/lib/hbase-server-0.
98.14-hadoop2.jar:/usr/local/hbase/lib/hbase-server-0.98.14-hadoop2-tests.jar:/usr/local/hbase/lib/hbase-shell-0.98.14-hadoop2.jar:/usr/l
ocal/hbase/lib/hbase-testing-util-0.98.14-hadoop2.jar:/usr/local/hbase/lib/hbase-thrift-0.98.14-hadoop2.jar:/usr/local/hbase/lib/high-sc
ale-lib-1.1.1.jar:/usr/local/hbase/lib/htrace-core-2.04.jar:/usr/local/hbase/lib/httpclient-4.1.3.jar:/usr/local/hbase/lib/httpcore-4.1.3.
jar:/usr/local/hbase/lib/jackson-core-asl-1.8.8.jar:/usr/local/hbase/lib/jackson-jaxrs-1.8.8.jar:/usr/local/hbase/lib/jackson-mapper-asl-
1.8.8.jar:/usr/local/hbase/lib/jackson-xc-1.8.8.jar:/usr/local/hbase/lib/jamon-runtime-2.3.1.jar:/usr/local/hbase/lib/jasper-compiler-5.5
.23.jar:/usr/local/hbase/lib/jasper-runtime-5.5.23.jar:/usr/local/hbase/lib/javax.inject-1.jar:/usr/local/hbase/lib/javax.servlet-3.1.jar
17/12/11 07:21:05 INFO node.PollingPropertiesFileConfigurationProvider: Configuration provider starting
17/12/11 07:21:05 INFO node.PollingPropertiesFileConfigurationProvider: Reloading configuration file:/home/acadgild/apache-flume-1.8.0-bi
n/conf/filecopy.conf
17/12/11 07:21:05 INFO conf.FlumeConfiguration: Processing:hdfsdest
17/12/11 07:21:05 INFO conf.FlumeConfiguration: Processing:hdfsdest
17/12/11 07:21:05 INFO conf.FlumeConfiguration: Processing:hdfsdest
17/12/11 07:21:05 INFO conf.FlumeConfiguration: Added sinks: hdfsdest Agent: agent1
17/12/11 07:21:05 INFO conf.FlumeConfiguration: Post-validation flume configuration contains configuration for agents: [agent1]
17/12/11 07:21:05 INFO node.AbstractConfigurationProvider: Creating channels
17/12/11 07:21:05 INFO channel.DefaultChannelFactory: Creating instance of channel mychannel type memory
17/12/11 07:21:05 INFO node.AbstractConfigurationProvider: Created channel mychannel
17/12/11 07:21:05 INFO source.DefaultSourceFactory: Creating instance of source mysrc, type exec
17/12/11 07:21:05 INFO sink.DefaultSinkFactory: Creating instance of sink: hdfsdest, type: hdfs
17/12/11 07:21:05 INFO node.AbstractConfigurationProvider: Channel mychannel connected to [mysrc, hdfsdest]
17/12/11 07:21:05 INFO node.Application: Starting new configuration:{ sourceRunners:[mysrc-EventDrivenSourceRunner: { source:org.apache.f
lume.source.ExecSource{name:mysrc,state:IDLE}}] sinkRunners:[hdfsdest-SinkRunner: { policy:org.apache.flume.sink.DefaultSinkProcessor@d5
4b076 counterGroup:{ name:null counters:{}} ] } } channels:[mychannel=org.apache.flume.channel.MemoryChannel{name: mychannel}] }
17/12/11 07:21:05 INFO node.Application: Starting channel mychannel
17/12/11 07:21:05 INFO instrumentation.MonitoredCounterGroup: Monitored counter group for type: CHANNEL, name: mychannel: Successfully re
gistered new MBean.
17/12/11 07:21:05 INFO instrumentation.MonitoredCounterGroup: Component type: CHANNEL, name: mychannel started
17/12/11 07:21:05 INFO node.Application: Starting Sink hdfsdest
17/12/11 07:21:05 INFO node.Application: Starting Source mysrc
17/12/11 07:21:05 INFO source.ExecSource: Exec source starting with command: hadoop dfs -put /home/acadgild/StatewiseDistrictwisePhysical
Progress.xml /flume_import
17/12/11 07:21:05 INFO instrumentation.MonitoredCounterGroup: Monitored counter group for type: SOURCE, name: mysrc: Successfully regist
ered new MBean.
17/12/11 07:21:05 INFO instrumentation.MonitoredCounterGroup: Component type: SOURCE, name: mysrc started
17/12/11 07:21:05 INFO instrumentation.MonitoredCounterGroup: Monitored counter group for type: SINK, name: hdfsdest: Successfully regist
ered new MBean.
17/12/11 07:21:05 INFO instrumentation.MonitoredCounterGroup: Component type: SINK, name: hdfsdest started
17/12/11 07:21:10 INFO source.ExecSource: Command [hadoop dfs -put /home/acadgild/StatewiseDistrictwisePhysicalProgress.xml /flume_import
] exited with 0
```

2. Verify that file is copied using the HDFS commands below

The command below confirms the directory is created

`hadoop fs -ls /flume_import`

```
[acadgild@localhost ~]$ hadoop fs -ls /flume_import
17/12/11 07:25:37 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
-rw-r--r-- 1 acadgild supergroup 717414 2017-12-11 07:21 /flume_import
```

3. Create folders in HDFS to store query outputs

`hadoop fs -mkdir districts_having_100percent_objectives`

`hadoop fs -mkdir districts_having_80percent_objectives`

```

[acadgild@localhost ~]$ hadoop fs -mkdir districts_having_100percent_objectives
17/12/11 07:40:14 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
[acadgild@localhost ~]$ hadoop fs -mkdir districts_having_80percent_objectives
17/12/11 07:40:23 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
[acadgild@localhost ~]$ hadoop fs -ls
17/12/11 07:40:28 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Found 9 items
drwxr-xr-x - acadgild supergroup 0 2015-11-20 11:46 Pictures
drwxr-xr-x - acadgild supergroup 0 2017-11-14 22:54 _sqoop
drwxr-xr-x - acadgild supergroup 0 2017-11-18 20:26 Company
drwxr-xr-x - acadgild supergroup 0 2017-12-11 07:40 districts_having_100percent_objectives
drwxr-xr-x - acadgild supergroup 0 2017-12-11 07:40 districts_having_80percent_objectives
drwxr-xr-x - acadgild supergroup 0 2017-12-11 07:37 folders
drwxr-xr-x - acadgild supergroup 0 2017-10-08 20:54 hadoop
drwxr-xr-x - acadgild supergroup 0 2015-11-17 02:03 oozie-acad
drwxr-xr-x - acadgild supergroup 0 2015-11-17 02:00 share
[acadgild@localhost ~]$

```

#### 4. Create mysql table to store the results of query

Start mysql using command

`sudo service mysqld start`

Login to mysql using

`mysql -u root`

Create Database bpl\_results

`create database bpl_results;`

`use bpl_results;`

Create tables `districts_having_100percent_objectives` and `districts_having_80percent_objectives` as below:

`create table districts_having_100percent_objectives`

`(`  
`name varchar(40)`

`);`

`create table districts_having_80percent_objectives`

`(`  
`name varchar(40)`

`);`

```
mysql> create database bpl_results;
Query OK, 1 row affected (0.00 sec)

mysql> use bpl_results;
Database changed
mysql> create table districts_having_100percent_objectives
-> (
->     name varchar(40)
-> );
Query OK, 0 rows affected (0.00 sec)

mysql> create table districts_having_80percent_objectives
-> (
->     name varchar(40)
-> );
Query OK, 0 rows affected (0.00 sec)

mysql> show tables;
+-----+
| Tables_in_bpl_results |
+-----+
| districts_having_100percent_objectives |
| districts_having_80percent_objectives |
+-----+
2 rows in set (0.00 sec)

mysql> █
```

##### 5. PIG query to process XML and store into PIG table

Load data from HDFS to PIG alias row\_physical\_progress using below query:

DEFINE XPath org.apache.pig.piggybank.evaluation.xml.XPath;

row\_physical\_progress = LOAD 'hdfs://localhost:9000/flume\_import' USING org.apache.pig.piggybank.storage.XMLLoader('row') as (row:chararray);

```
grunt> DEFINE XPath org.apache.pig.piggybank.evaluation.xml.XPath;
grunt> row_physical_progress = LOAD 'hdfs://localhost:9000/flume_import' USING org.apache.pig.piggybank.storage.XMLLoader('row') as (row:chararray);
2017-12-11 08:14:06,102 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker.persist.jobstatus.hours is deprecated. Instead, use mapreduce.jobtracker.persist.jobstatus.hours
2017-12-11 08:14:06,102 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.heartbeats.in.second is deprecated. Instead, use mapreduce.jobtracker.heartbeats.in.second
2017-12-11 08:14:06,102 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - jobclient.completion.poll.interval is deprecated. Instead, use mapreduce.client.completion.pollinterval
2017-12-11 08:14:06,102 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.tasktracker.tasks.sleep-time-before-sigkill is deprecated. Instead, use mapreduce.tasktracker.tasks.sleep-time-before-sigkill
2017-12-11 08:14:06,102 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker.http.address is deprecated. Instead, use mapreduce.jobtracker.http.address
```

Next, iterate over each row and load into alias physical\_progress which has schema fields same as XML schema hyphen(-) are replaced with underscore (\_)

physical\_progress = FOREACH row\_physical\_progress GENERATE XPath(row, 'row/State\_Name') AS State\_Name,

XPath(row, 'row/District\_Name') AS District\_Name,

XPath(row, 'row/Project\_Objectives\_IHHL\_BPL') AS

Project\_Objectives\_IHHL\_BPL,

XPath(row, 'row/Project\_Objectives\_IHHL\_APL') AS

Project\_Objectives\_IHHL\_APL,

XPath(row, 'row/Project\_Objectives\_IHHL\_TOTAL') AS

Project\_Objectives\_IHHL\_TOTAL,

XPath(row, 'row/Project\_Objectives\_SCW') AS Project\_Objectives\_SCW,

XPath(row, 'row/Project\_Objectives\_Anganwadi\_Toilets') AS

Project\_Objectives\_Anganwadi\_Toilets,

```

        XPath(row, 'row/Project_Objectives_RSM') AS Project_Objectives_RSM,
        XPath(row, 'row/Project_Objectives_PC') AS Project_Objectives_PC,
        XPath(row, 'row/Project_Performance-IHHL_BPL') AS
Project_Performance_IHHL_BPL,
        XPath(row, 'row/Project_Performance-IHHL_APL') AS
Project_Performance_IHHL_APL,
        XPath(row, 'row/Project_Performance-IHHL_TOTAL') AS
Project_Performance_IHHL_TOTAL,
        XPath(row, 'row/Project_Performance-SCW') AS Project_Performance_SCW,
        XPath(row, 'row/Project_Performance-School_Toilets') AS
Project_Performance_School_Toilets,
        XPath(row, 'row/Project_Performance-Anganwadi_Toilets') AS
Project_Performance_Anganwadi_Toilets,
        XPath(row, 'row/Project_Performance-RSM') AS Project_Performance_RSM,
        XPath(row, 'row/Project_Performance-PC') AS Project_Performance_PC;

```

```

grunt> physical_progress = FOREACH row_physical_progress GENERATE XPath(row, 'row/State_Name') AS State_Name,
>> XPath(row, 'row/District_Name') AS District_Name,
>> XPath(row, 'row/Project_Objectives_IHHL_BPL') AS Project_Objectives_IHHL_BPL,
>> XPath(row, 'row/Project_Objectives_IHHL_APL') AS Project_Objectives_IHHL_APL,
>> XPath(row, 'row/Project_Objectives_IHHL_TOTAL') AS Project_Objectives_IHHL_TOTAL,
>> XPath(row, 'row/Project_Objectives_SCW') AS Project_Objectives_SCW,
>>
>> XPath(row, 'row/Project_Objectives_Anganwadi_Toilets') AS Project_Objectives_Anganwadi_Toilets,
>> XPath(row, 'row/Project_Objectives_RSM') AS Project_Objectives_RSM,
>> XPath(row, 'row/Project_Objectives_PC') AS Project_Objectives_PC,
>> XPath(row, 'row/Project_Performance-IHHL_BPL') AS Project_Performance_IHHL_BPL,
>> XPath(row, 'row/Project_Performance-IHHL_APL') AS Project_Performance_IHHL_APL,
>> XPath(row, 'row/Project_Performance-IHHL_TOTAL') AS Project_Performance_IHHL_TOTAL,
>> XPath(row, 'row/Project_Performance-SCW') AS Project_Performance_SCW,
>> XPath(row, 'row/Project_Performance-School_Toilets') AS Project_Performance_School_Toilets,
>> XPath(row, 'row/Project_Performance-Anganwadi_Toilets') AS Project_Performance_Anganwadi_Toilets,
>> XPath(row, 'row/Project_Performance-RSM') AS Project_Performance_RSM,
>> XPath(row, 'row/Project_Performance-PC') AS Project_Performance_PC;
grunt> █

```

#### 6. PIG Query to find out the districts who achieved 100 percent objective in BPL cards

Here first filter the records where Project\_Objectives\_IHHL\_BPL is equal to Project\_Performance\_IHHL\_BPL

```
physical_progress_100_percent_bpl = FILTER physical_progress BY
Project_Objectives_IHHL_BPL == Project_Performance_IHHL_BPL;
```

Next, Select only District\_Name field using command below:

```
district_100_percent_bpl = FOREACH physical_progress_100_percent_bpl GENERATE
District_Name;
```

Next, Store into HDFS directory districts\_having\_100percent\_objectives using command below:

```
STORE district_100_percent_bpl INTO
'hdfs://localhost:9000/districts_having_100percent_objectives'
```

```

grunt> physical_progress_100_percent_bpl = FILTER physical_progress BY Project_Objectives_IHHL_BPL == Project_Performance_IHHL_BPL;
grunt> district_100_percent_bpl = FOREACH physical_progress_100_percent_bpl GENERATE District_Name;
grunt> STORE district_100_percent_bpl INTO 'hdfs://localhost:9000/districts_having_100percent_objectives'
>> ;
2017-12-11 08:21:01,669 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapreduce.job.counters.limit is deprecated. Instead, use mapreduce.job.counters.max
2017-12-11 08:21:01,670 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2017-12-11 08:21:01,670 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2017-12-11 08:21:01,752 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.textoutputformat.separator is deprecated. Instead, use mapreduce.output.textoutputformat.separator
2017-12-11 08:21:01,822 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: FILTER
2017-12-11 08:21:01,933 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2017-12-11 08:21:01,937 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS

```

#### < ----- Intermediate Logs ----- >

```

HadoopVersion 2.2.0 PigVersion 0.14.0
Pig acadgild
User 2017-12-11 08:21:02
StartedAt 2017-12-11 08:21:43
FinishedAt 2017-12-11 08:21:43
Features FILTER
Success!

Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime M
edianReduceTime Alias Feature Outputs
job_local1176682894_0001 1 0 n/a n/a n/a n/a 0 0 0 district_100_percent_bpl,
physical_progress,physical_progress_100_percent_bpl,row_physical_progress
0percent_objectives,
Input(s):
Successfully read 607 records (729702 bytes) from: "hdfs://localhost:9000/flume_import"
Output(s):
Successfully stored 70 records (686 bytes) in: "hdfs://localhost:9000/districts_having_100percent_objectives"
Counters:
Total records written : 70
Total bytes written : 686
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0
Job DAG:
job_local1176682894_0001
2017-12-11 08:21:43,547 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize JVM Metrics with processName=JobTracker
, sessionId= - already initialized
2017-12-11 08:21:43,554 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize JVM Metrics with processName=JobTracker
, sessionId= - already initialized
2017-12-11 08:21:43,554 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize JVM Metrics with processName=JobTracker
, sessionId= - already initialized
2017-12-11 08:21:43,584 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
grunt>

```

#### 7. Verify that results are stored in HDFS

The following command shows that folders are created under  
districts\_having\_100percent\_objectives  
hadoop fs -ls /districts\_having\_100percent\_objectives

Next, use the following HDFS command to show the results

hadoop fs -ls /districts\_having\_100percent\_objectives/part-m-00000

```

[acadgild@localhost ~]$ hadoop fs -ls /districts_having_100percent_objectives
17/12/11 08:29:01 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Found 2 items
-rw-r--r-- 3 acadgild supergroup 0 2017-12-11 08:21 /districts_having_100percent_objectives/_SUCCESS
-rw-r--r-- 3 acadgild supergroup 686 2017-12-11 08:21 /districts_having_100percent_objectives/part-m-00000 ←

```



```
[acadgild@localhost ~]$ hadoop fs -cat /districts_having_100percent_objectives/part-m-00000
17/12/11 08:32:05 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform.
applicable
NIZAMABAD
TIRAP
HAILAKANDI
MADHUBANI
NORTH GOA
AHMEDABAD
DANGS
NAVSARI
PORBANDAR
SURAT
FARIDABAD
HISAR
JHAJJAR
MAHENDRAGARH
PANCHKULA
PANIPAT
ROHTAK
SIRSA
HAMIRPUR
KINNAUR
KULLU
LAHAUL & SPITI
SHIMLA
SOLAN
UNA
DEOGHAR
LOHARDAGA
HASSAN
MANGALORE (DAKSHINA KANNADA)
UDUPI
ALAPPUZHA
KOLLAM
KOTTAYAM
KOZHIKODE
PALAKKAD
PATHANAMTHITTA
WAYANAD
GADCHIROLI
SINDHUDURG
WEST GARO HILLS
CHAMPHAI
LAWNGTLAI
HANUMANGARH
ERODE
KARUR
NAMAKKAL
TIRUCHIRAPPALLI
TIRUVANNAMALAI
DHALAI
SOUTH TRIPURA
WEST TRIPURA
AMBEDKAR NAGAR
BALRAMPUR
BAREILLY
BIJNOR
BUDAUN
ETAWAH
FARRUKHABAD
FIROZABAD
GHAZIABAD
HARDOI
JYOTIBA PHULE NAGAR
LUCKNOW
MAHARAJGANJ
MAHOBA
MORADABAD
MUZAFFARNAGAR
PILIBHIT
SONBHADRA
SULTANPUR
[acadgild@localhost ~]$
```

8. Use sqoop command to export data from HDFS into mysql table districts\_having\_100percent\_objectives in database bpl\_results

The following sqoop command is used to export data from HDFS folder districts\_having\_100percent\_objectives into already created mysql table 'districts\_having\_100percent\_objectives'

```

[acagild@localhost ~]$ sqoop export --connect jdbc:mysql://localhost/bpl_results --username 'root' --table 'districts_having_100percent_objectives' --export-dir '/districts_having_100percent_objectives' --input-fields-terminated-by ',' -m 1 --columns name;
Warning: /usr/local/sqoop/./hcatalog does not exist! HCatalog jobs will fail.
Please set $HCAT_HOME to the root of your HCatalog installation.
Warning: /usr/local/sqoop/./accumulo does not exist! Accumulo imports will fail.
Please set $ACCUMULO_HOME to the root of your Accumulo installation.
Warning: /usr/local/sqoop/./zookeeper does not exist! Accumulo imports will fail.
Please set $ZOOKEEPER_HOME to the root of your Zookeeper installation.
2017-12-11 08:37:30,222 INFO [main] sqoop.Sqoop: Running Sqoop version: 1.4.6
2017-12-11 08:37:30,638 INFO [main] manager.MySQLManager: Preparing to use a MySQL streaming resultset.
2017-12-11 08:37:30,638 INFO [main] tool.CodeGenTool: Beginning code generation
2017-12-11 08:37:31,334 INFO [main] manager.SqlManager: Executing SQL statement: SELECT t.* FROM `districts_having_100percent_objectives` AS t LIMIT 1
2017-12-11 08:37:31,399 INFO [main] manager.SqlManager: Executing SQL statement: SELECT t.* FROM `districts_having_100percent_objectives` AS t LIMIT 1
2017-12-11 08:37:31,409 INFO [main] orm.CompilationManager: HADOOP_MAPRED_HOME is /usr/local/hadoop-2.6.0
Note: /tmp/sqoop-acagild/compile/108fa63cce4b2e2786d49ebff92261ec/districts_having_100percent_objectives.java uses or overrides a deprecated API.
Note: Recompile with -Xlint:deprecation for details.
2017-12-11 08:37:38,448 INFO [main] orm.CompilationManager: Writing jar file: /tmp/sqoop-acagild/compile/108fa63cce4b2e2786d49ebff92261ec/districts_having_100percent_objectives.jar
2017-12-11 08:37:38,505 INFO [main] mapreduce.ExportJobBase: Beginning export of districts_having_100percent_objectives
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/local/hbase/lib/slf4j-log4j12-1.6.4.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/local/hadoop-2.6.0/share/hadoop/common/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
2017-12-11 08:37:39,071 WARN [main] util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable

```

#### <----- Intermediate Logs ----->

```

2017-12-11 08:37:45,312 INFO [main] Configuration.deprecation: user.name is deprecated. Instead, use mapreduce.job.user.name
2017-12-11 08:37:45,312 INFO [main] Configuration.deprecation: mapred.reduce.tasks is deprecated. Instead, use mapreduce.job.reduces
2017-12-11 08:37:45,312 INFO [main] Configuration.deprecation: mapred.cache.files.filesizes is deprecated. Instead, use mapreduce.job.cache.files.filesizes
2017-12-11 08:37:45,555 INFO [main] mapreduce.JobSubmitter: Submitting tokens for job: job_1512956846592_0001
2017-12-11 08:37:46,561 INFO [main] impl.YarnClientImpl: Submitted application application_1512956846592_0001 to ResourceManager at /0.0.0:8032
2017-12-11 08:37:46,623 INFO [main] mapreduce.Job: The url to track the job: http://http://localhost:8088/proxy/application_1512956846592_0001/
2017-12-11 08:37:46,623 INFO [main] mapreduce.Job: Running job: job_1512956846592_0001
2017-12-11 08:38:06,043 INFO [main] mapreduce.Job: Job job_1512956846592_0001 running in uber mode : false
2017-12-11 08:38:06,046 INFO [main] mapreduce.Job: map 0% reduce 0%
2017-12-11 08:38:16,007 INFO [main] mapreduce.Job: map 100% reduce 0%
2017-12-11 08:38:17,036 INFO [main] mapreduce.Job: Job job_1512956846592_0001 completed successfully

```

## 9. Verify Result in Mysql

Use the following command in mysql to verify results in mysql

```
select * from districts_having_100percent_objectives;
```



```
mysql> select * from districts_having_100percent_objectives;
```

name
NIZAMABAD
TIRAP
HAILAKANDI
MADHUBANI
NORTH GOA
AHMEDABAD
DANGS
NAVSARI
PORBANDAR
SURAT
FARIDABAD
HISAR
JHAJJAR
MAHENDRAGARH
PANCHKULA
PANIPAT
ROHTAK
SIRSA
HAMIRPUR
KINNAUR
KULLU
LAHAUL & SPITI
SHIMLA
SOLAN
UNA
DEOGHAR
LOHARDAGA
HASSAN
MANGALORE (DAKSHINA KANNADA)
UDUPI
ALAPPUZHA
KOLLAM
KOTTAYAM
KOZHIKODE

```

| PALAKKAD
| PATHANAMTHITTA
| WAYANAD
| GADCHIROLI
| SINDHUDURG
| WEST GARO HILLS
| CHAMPHAI
| LAWNGTLAI
| HANUMANGARH
| ERODE
| KARUR
| NAMAKKAL
| TIRUCHIRAPPALLI
| TIRUVANNAMALAI
| DHALAI
| SOUTH TRIPURA
| WEST TRIPURA
| AMBEDKAR NAGAR
| BALRAMPUR
| BAREILLY
| BIJNOR
| BUDAUN
| ETAWAH
| FARRUKHABAD
| FIROZABAD
| GHAZIABAD
| HARDOI
| JYOTIBA PHULE NAGAR
| LUCKNOW
| MAHARAJGANJ
| MAHOBA
| MORADABAD
| MUZAFFARNAGAR
| PILIBHIT
| SONBHADRA
| SULTANPUR
+-----+
70 rows in set (0.00 sec)

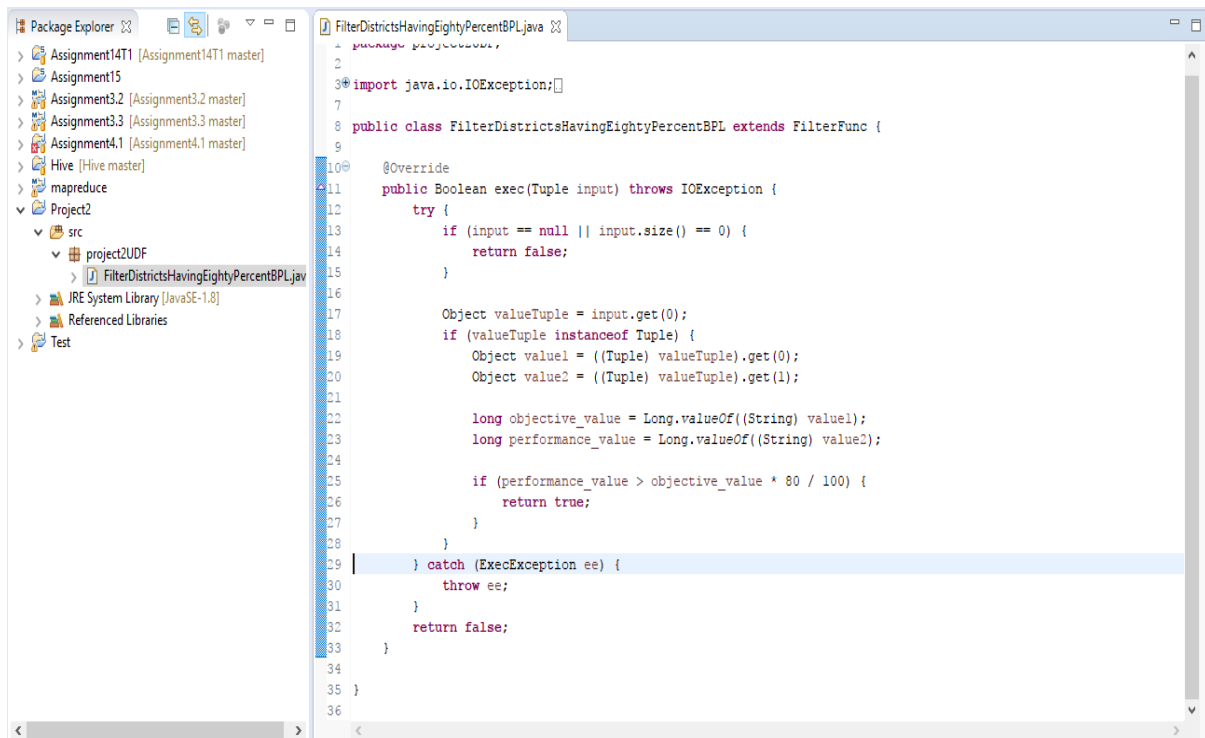
```

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

**Export the results to MySQL using Sqoop.**

1. Create a Java project Project2 and Write a Java class FilterDistrictsHavingEightyPercentBPL in eclipse which will filter those tuples for which 80 percent objective in BPL cards are achieved. The logic put in exec method is value of Project\_Performance\_IHHL\_BPL is equal to more than 80% of Project\_Objectives\_IHHL\_BPL.

Export the project to Project2.jar



2. Write PIG query to find out the districts who achieved 80 percent objective in BPL cards  
Register the Jar Project2.jar for the UDF created in step11

REGISTER /home/acadgild/pig/Project2.jar;

Next, using the UDF filter those tuple for which Project\_Performance\_IHHL\_BPL is equal to more than 80% of Project\_Objectives\_IHHL\_BPL

physical\_progress\_80\_percent\_bpl = FILTER physical\_progress BY  
project2UDF.FilterDistrictsHavingEightyPercentBPL(TOTUPLE(Project\_Objectives\_IHHL\_BPL,  
Project\_Performance\_IHHL\_BPL));

Next, Select only District\_Name field using command below:

district\_80\_percent\_bpl = FOREACH physical\_progress\_80\_percent\_bpl GENERATE  
District\_Name;

Next, Store into HDFS directory districts\_having\_100percent\_objectives using command below:  
STORE district\_80\_percent\_bpl INTO  
'hdfs://localhost:9000/districts\_having\_80percent\_objectives';

```

grunt> REGISTER /home/acadgild/pig/Project2.jar;
2017-12-11 09:43:53,699 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use
dfs.bytes-per-checksum
2017-12-11 09:43:53,699 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.de
faultFS
2017-12-11 09:43:53,699 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapreduce.job.counters.limit is deprecated. Inste
ad, use mapreduce.job.counters.max
2017-12-11 09:43:53,769 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapreduce.job.counters.limit is deprecated. Inste
ad, use mapreduce.job.counters.max
2017-12-11 09:43:53,769 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use
dfs.bytes-per-checksum
2017-12-11 09:43:53,769 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.de
faultFS
grunt> physical_progress_80_percent_bpl = FILTER physical_progress BY project2UDF.FilterDistrictsHavingEightyPercentBPL(TOTUPLE(Project_0
bjectives_IHHL_BPL, Project_Performance_IHHL_BPL));
grunt> district_80_percent_bpl = FOREACH physical_progress_80_percent_bpl GENERATE District Name;
grunt> STORE district_80_percent_bpl INTO 'hdfs://localhost:9000/districts_having_80percent_objectives';
2017-12-11 09:45:18,472 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapreduce.job.counters.limit is deprecated. Inste
ad, use mapreduce.job.counters.max
2017-12-11 09:45:18,472 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use
dfs.bytes-per-checksum
2017-12-11 09:45:18,472 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.de
faultFS
2017-12-11 09:45:18,572 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: FILTER
2017-12-11 09:45:18,650 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use
dfs.bytes-per-checksum
2017-12-11 09:45:18,650 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.de
faultFS
2017-12-11 09:45:18,653 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapreduce.job.counters.limit is deprecated. Inste
ad, use mapreduce.job.counters.max

```

### <----- Intermediate Logs ----->

```

HadoopVersion  PigVersion  UserId  StartedAt  FinishedAt  Features
2.2.0  0.14.0  acadgild  2017-12-11 09:45:18  2017-12-11 09:46:02  FILTER

Success!

Job Stats (time in seconds):
JobId  Maps  Reduces  MaxMapTime  MinMapTime  AvgMapTime  MedianMapTime  MaxReduceTime  MinReduceTime  AvgReduceTime  M
edianReduceTime  Alias  Feature  Outputs
job_local427933718_0002  1  0  n/a  n/a  n/a  n/a  0  0  0  district_80_percent_bpl,physical_
progress,physical_progress_80_percent_bpl,row_physical_progress  MAP_ONLY  hdfs://localhost:9000/districts_having_80percent_objectiv
es,

Input(s):
Successfully read 607 records (1447116 bytes) from: "hdfs://localhost:9000/flume_import"

Output(s):
Successfully stored 349 records (4038 bytes) in: "hdfs://localhost:9000/districts_having_80percent_objectives"

Counters:
Total records written : 349
Total bytes written : 4038
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_local427933718_0002

2017-12-11 09:46:02,451 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize JVM Metrics with processName=JobTracker
, sessionId= - already initialized
2017-12-11 09:46:02,452 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize JVM Metrics with processName=JobTracker
, sessionId= - already initialized
2017-12-11 09:46:02,455 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize JVM Metrics with processName=JobTracker
, sessionId= - already initialized
2017-12-11 09:46:02,474 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
grunt>

```

### 3. Verify that results are stored in HDFS

The following command shows that folders are created under  
districts\_having\_100percent\_objectives

hadoop fs -ls /districts\_having\_80percent\_objectives

Next, use the following HDFS command to show the results

hadoop fs -ls /districts\_having\_80percent\_objectives/part-m-00000

```
[acadgild@localhost ~]$ hadoop fs -ls /districts_having_80percent_objectives
17/12/11 09:53:08 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Found 2 items
-rw-r--r-- 3 acadgild supergroup 0 2017-12-11 09:46 /districts_having_80percent_objectives/_SUCCESS
-rw-r--r-- 3 acadgild supergroup 3352 2017-12-11 09:46 /districts_having_80percent_objectives/part-m-00000
[acadgild@localhost ~]$ hadoop fs -cat /districts_having_80percent_objectives/part-m-00000
17/12/11 09:53:24 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
ANANTAPUR
CHITTOOR
CUDDAPAH
EAST GODAVARI
KARIMNAGAR
KHAMMAM
KRISHNA
KURNOOL
MEDAK
NALGONDA
NIZAMABAD
RANGAREDDI
WARANGAL
WEST GODAVARI
DIBANG VALLEY
LOHIT
TIRAP
BAGSHA
CACHAR
DIBRUGARH
GOALPARA
GOLAGHAT
HAILAKANDI
```

4. Use sqoop command to export data from HDFS into mysql table districts\_having\_80percent\_objectives in database bpl\_results

The following sqoop command is used to export data from HDFS folder districts\_having\_80percent\_objectives into already created mysql table 'districts\_having\_80percent\_objectives'

Screenshots are as below:

sqoop export --connect jdbc:mysql://localhost/bpl\_results --username 'root' --table 'districts\_having\_80percent\_objectives' --export-dir '/districts\_having\_80percent\_objectives' --input-fields-terminated-by ';' -m 1 --columns name

```
[acadgild@localhost ~]$ sqoop export --connect jdbc:mysql://localhost/bpl_results --username 'root' --table 'districts_having_80percent_o
bjectives' --export-dir '/districts_having_80percent_objectives' --input-fields-terminated-by ';' -m 1 --columns name
Warning: /usr/local/sqoop/./hcatalog does not exist! HCatalog jobs will fail.
Please set $HCAT_HOME to the root of your HCatalog installation.
Warning: /usr/local/sqoop/./accumulo does not exist! Accumulo imports will fail.
Please set $ACCUMULO_HOME to the root of your Accumulo installation.
Warning: /usr/local/sqoop/./zookeeper does not exist! Accumulo imports will fail.
Please set $ZOOKEEPER_HOME to the root of your Zookeeper installation.
2017-12-11 09:57:19,250 INFO [main] sqoop.Sqoop: Running Sqoop version: 1.4.6
2017-12-11 09:57:20,156 INFO [main] manager.MySQLManager: Preparing to use a MySQL streaming resultset.
2017-12-11 09:57:20,156 INFO [main] tool.CodeGenTool: Beginning code generation
2017-12-11 09:57:21,045 INFO [main] manager.SqlManager: Executing SQL statement: SELECT t.* FROM `districts_having_80percent_objectives`
AS t LIMIT 1
2017-12-11 09:57:21,145 INFO [main] manager.SqlManager: Executing SQL statement: SELECT t.* FROM `districts_having_80percent_objectives`
AS t LIMIT 1
2017-12-11 09:57:21,159 INFO [main] orm.CompilationManager: HADOOP_MAPRED_HOME is /usr/local/hadoop-2.6.0
Note: /tmp/sqoop-acadgild/compile/297be123c7d8b9a6a6fefce4bef412e3/districts_having_80percent_objectives.java uses or overrides a depreca
ted API.
Note: Recompile with -Xlint:deprecation for details.
2017-12-11 09:57:25,563 INFO [main] orm.CompilationManager: Writing jar file: /tmp/sqoop-acadgild/compile/297be123c7d8b9a6a6fefce4bef412
e3/districts_having_80percent_objectives.jar
2017-12-11 09:57:25,608 INFO [main] mapreduce.ExportJobBase: Beginning export of districts_having_80percent_objectives
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/local/hbase/lib/slf4j-log4j12-1.6.4.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/local/hadoop-2.6.0/share/hadoop/common/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBin
der.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
2017-12-11 09:57:26,407 WARN [main] util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java
classes where applicable
2017-12-11 09:57:26,418 INFO [main] Configuration.deprecation: mapred.jar is deprecated. Instead, use mapreduce.job.jar
2017-12-11 09:57:28,275 INFO [main] Configuration.deprecation: mapred.reduce.tasks.speculative.execution is deprecated. Instead, use map
reduce.reduce.speculative
2017-12-11 09:57:28,298 INFO [main] Configuration.deprecation: mapred.map.tasks.speculative.execution is deprecated. Instead, use mapred
uce.map.speculative
2017-12-11 09:57:28,304 INFO [main] Configuration.deprecation: mapred.map.tasks is deprecated. Instead, use mapreduce.job.maps
```

5. Verify Result in Mysql

Use the following command in mysql to verify results in mysql

```
select * from districts_having_80percent_objectives
```

```
mysql> select * from districts_having_80percent_objectives;
```

name
ANANTAPUR
CHITTOOR
CUDDAPAH
EAST GODAVARI
KARIMNAGAR
KHAMMAM
KRISHNA
KURNOOL
MEDAK
NALGONDA
NIZAMABAD
RANGAREDDI
WARANGAL
WEST GODAVARI
DIBANG VALLEY
LOHIT
TIRAP
BAGSHA
CACHAR
DIBRUGARH
GOALPARA
GOLAGHAT
HAILAKANDI
JORHAT
KAMRUP
KARIMGANJ
KOKRAJHAR
LAKHIMPUR
MARIGAON
NAGAON
SIBSAGAR
SONITPUR
TINSUKIA
BEGUSARAI



MUZAFFARPUR	
SAHARSA	
VAISHALI	
DHAMTARI	
JASHPUR	
KANKER	
KORBA	
KORIYA	
SURGUJA	
NORTH GOA	
AHMEDABAD	
AMRELI	
ANAND	
BANAS KANTHA	
BHARUCH	
BHAVNAGAR	
DAHOD	
DANGS	
GANDHINAGAR	
JAMNAGAR	
JUNAGADH	
KACHCHH	
KHEDA	
MAHESANA	
NARMADA	
NAVSARI	
PANCH MAHALS	
PATAN	
PORBANDAR	
RAJKOT	
SABAR KANTHA	
SURAT	
SURENDRANAGAR	
VADODARA	
VALSAD	
AMBALA	
BHIWANI	
FARIDABAD	

: :

```
| PILIBHIT  
| PRATAPGARH  
| RAE BARELI  
| RAMPUR  
| SAHARANPUR  
| SANT RAVIDAS NAGAR( BHADOHI)  
| SHAHJAHANPUR  
| SHRAVASTI  
| SIDDHARTHANAGAR  
| SITAPUR  
| SONBHADRA  
| SULTANPUR  
| UNNAO  
| VARANASI  
| BAGESHWAR  
| CHAMOLI  
| DEHRADUN  
| HARIDWAR  
| NAINITAL  
| PITHORAGARH  
| RUDRAPRAYAG  
| TEHRI GARHWAL  
| UDHAM SINGH NAGAR  
| UTTARKASHI  
| BARDHAMAN  
| DAKSHIN DINAJPUR  
| HOOGHLY  
| HOWRAH  
| JALPAIGURI  
| MIDNAPUR EAST  
| MIDNAPUR WEST  
| NADIA  
| NORTH 24 PARAGANAS  
| SOUTH 24 PARAGANAS  
+-----+  
349 rows in set (0.00 sec)  
  
mysql> █
```