

**Exp.No.: 4****Create UDF in PIG****Step-by-step installation of Apache Pig on Hadoop cluster on Ubuntu Pre-requisite:**

- Ubuntu 16.04 or higher version running (I have installed Ubuntu on Oracle VM (Virtual Machine) VirtualBox),
- Run Hadoop on ubuntu (I have installed Hadoop 3.2.1 on Ubuntu 16.04). You may refer to my blog “How to install Hadoop installation” click [here](#) for Hadoop installation).

**Pig installation steps****Step 1: Login into Ubuntu**

```
hadoop@hadoop-VirtualBox:~$ $ wget https://dlcdn.apache.org/pig/pig-0.16.0/pig-0.16.0.tar.gz
$: command not found
hadoop@hadoop-VirtualBox:~$ wget https://dlcdn.apache.org/pig/pig-0.16.0/pig-0.16.0.tar.gz
--2022-06-21 11:57:52-- https://dlcdn.apache.org/pig/pig-0.16.0/pig-0.16.0.tar.gz
Resolving dlcdn.apache.org (dlcdn.apache.org)... 151.101.2.132, 2a04:4e42::644
Connecting to dlcdn.apache.org (dlcdn.apache.org)|151.101.2.132|:443... connecte
d.
HTTP request sent, awaiting response... 200 OK
Length: 177279333 (169M) [application/x-gzip]
Saving to: 'pig-0.16.0.tar.gz.1'

pig-0.16.0.tar.gz.1  94%[=====> ] 158.94M  5.19MB/s  eta 2s
```

**Step 2:** Go to <https://pig.apache.org/releases.html> and copy the path of the latest version of pig that you want to install. Run the following command to download Apache Pig in Ubuntu:

```
$ wget https://dlcdn.apache.org/pig/pig-0.16.0/pig-0.16.0.tar.gz
```

**Step 3:** To untar pig-0.16.0.tar.gz file run the following command:

```
$ tar xvzf pig-0.16.0.tar.gz
```

**Step 4:** To create a pig folder and move pig-0.16.0 to the pig folder, execute the following command:

```
$ sudo mv /home/hadoop/pig-0.16.0 /home/hadoop/pig
```

**Step 5:** Now open the .bashrc file to edit the path and variables/settings for pig. Run the following command:

```
$ sudo nano .bashrc
```

Add the below given to .bashrc file at the end and save the file.

```
#PIG settings
export PIG_HOME=/home/hadoop/pig
export PATH=$PATH:$PIG_HOME/bin
```

```

PIG_CLASSPATH=$PIG_HOME/conf:$HADOOP_INSTALL/etc/hadoop/export
PIG_CONF_DIR=$PIG_HOME/confexport JAVA_HOME=/usr/lib/jvm/java-8-
openjdk-amd64export PIG_CLASSPATH=$PIG_CONF_DIR:$PATH#PIG setting ends

```

```

vinisha@ubuntu: ~
GNU nano 6.2 .bashrc
export HADOOP_OPTS="-Djava.library.path=$HADOOP_HOME/lib/native"
export HADOOP_STREAMING=$HADOOP_HOME/share/hadoop/tools/lib/hadoop-streaming-3.2.3.jar
export HADOOP_LOG_DIR=$HADOOP_HOME/logs
export PDSH_RCMD_TYPE=ssh

PIG settings
export PIG_HOME=/pig/pig-0.16.0
export PATH=$PATH:$PIG_HOME/bin
export PIG_CLASSPATH=$PIG_HOME/conf:$HADOOP_HOME/etc/hadoop
export PIG_CONF_DIR=$PIG_HOME/conf
export JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64
export PIG_CLASSPATH=$PIG_CONF_DIR:$PATH
PIG settings end

export HIVE_HOME=/home/vinisha/apache-hive-3.1.2-bin
export PATH=$PATH:$HIVE_HOME/bin

```

**Step 6:** Run the following command to make the changes effective in the .bashrc file:

```
$ source .bashrc
```

**Step 7:** To start all Hadoop daemons, navigate to the hadoop-3.2.1/sbin folder and run the following commands:

```
$ ./start-dfs.sh$ ./start-yarn$ jps
```

```

vinisha@ubuntu: $ sudo nano .bashrc
[sudo] password for vinisha:
vinisha@ubuntu: $ source ~/.bashrc
vinisha@ubuntu: $

2024-09-19 20:27:55,339 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCESS, Redirecting to Job History Server
2024-09-19 20:27:56,244 [main] INFO org.apache.hadoop.lpc.Client - Retrying connect to server: 0.0.0.0/0.0.0.0:10020. Already tried 0 time(s); retry policy is RetryUpToMaximumCountWithFixedSleep(maxRetries=10, sleepTime=1000 MILLISECONDS)
2024-09-19 20:27:57,357 [main] INFO org.apache.hadoop.lpc.Client - Retrying connect to server: 0.0.0.0/0.0.0.0:10020. Already tried 1 time(s); retry policy is RetryUpToMaximumCountWithFixedSleep(maxRetries=10, sleepTime=1000 MILLISECONDS)
2024-09-19 20:27:58,461 [main] INFO org.apache.hadoop.lpc.Client - Retrying connect to server: 0.0.0.0/0.0.0.0:10020. Already tried 2 time(s); retry policy is RetryUpToMaximumCountWithFixedSleep(maxRetries=10, sleepTime=1000 MILLISECONDS)
2024-09-19 20:27:59,488 [main] INFO org.apache.hadoop.lpc.Client - Retrying connect to server: 0.0.0.0/0.0.0.0:10020. Already tried 3 time(s); retry policy is RetryUpToMaximumCountWithFixedSleep(maxRetries=10, sleepTime=1000 MILLISECONDS)
2024-09-19 20:28:00,412 [main] INFO org.apache.hadoop.lpc.Client - Retrying connect to server: 0.0.0.0/0.0.0.0:10020. Already tried 4 time(s); retry policy is RetryUpToMaximumCountWithFixedSleep(maxRetries=10, sleepTime=1000 MILLISECONDS)
2024-09-19 20:28:01,417 [main] INFO org.apache.hadoop.lpc.Client - Retrying connect to server: 0.0.0.0/0.0.0.0:10020. Already tried 5 time(s); retry policy is RetryUpToMaximumCountWithFixedSleep(maxRetries=10, sleepTime=1000 MILLISECONDS)
2024-09-19 20:28:02,419 [main] INFO org.apache.hadoop.lpc.Client - Retrying connect to server: 0.0.0.0/0.0.0.0:10020. Already tried 6 time(s); retry policy is RetryUpToMaximumCountWithFixedSleep(maxRetries=10, sleepTime=1000 MILLISECONDS)
2024-09-19 20:28:03,430 [main] INFO org.apache.hadoop.lpc.Client - Retrying connect to server: 0.0.0.0/0.0.0.0:10020. Already tried 7 time(s); retry policy is RetryUpToMaximumCountWithFixedSleep(maxRetries=10, sleepTime=1000 MILLISECONDS)
2024-09-19 20:28:04,439 [main] INFO org.apache.hadoop.lpc.Client - Retrying connect to server: 0.0.0.0/0.0.0.0:10020. Already tried 8 time(s); retry policy is RetryUpToMaximumCountWithFixedSleep(maxRetries=10, sleepTime=1000 MILLISECONDS)
2024-09-19 20:28:05,445 [main] INFO org.apache.hadoop.lpc.Client - Retrying connect to server: 0.0.0.0/0.0.0.0:10020. Already tried 9 time(s); retry policy is RetryUpToMaximumCountWithFixedSleep(maxRetries=10, sleepTime=1000 MILLISECONDS)
2024-09-19 20:28:05,502 [main] WARN org.apache.hadoop.mapred.executionengine.nagreduceLayer.NagReduceLauncher - unable to retrieve Job to compute waiting deprecation.
2024-09-19 20:28:05,563 [main] INFO org.apache.hadoop.mapred.executionengine.nagreduceLayer.NagReduceLauncher - Success
2024-09-19 20:28:05,578 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publisher.enabled
2024-09-19 20:28:05,688 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2024-09-19 20:28:05,692 [main] INFO org.apache.hadoop.data.SchemaTableBackend - Key [pig.schematuple] was not set... will not generate code.
2024-09-19 20:28:05,922 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input files to process : 1
2024-09-19 20:28:05,923 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
(1,Job)
(2,Job)
(3,Job)
(4,Job)
grant>

```

**Step 8:** Now you can launch pig by executing the following command: \$ pi

**Step 9:** Now you are in pig and can perform your desired tasks on pig. You can come out of the pig by the quit command:

```
> quit;
```

## CREATE USER DEFINED FUNCTION

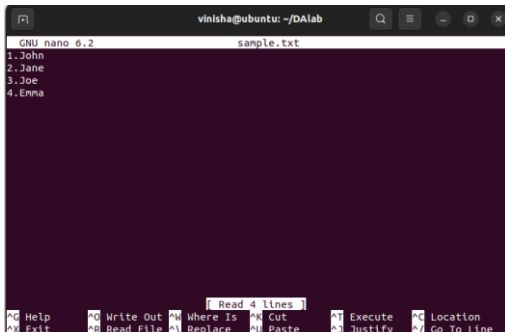
To create User Define Function in Apache Pig and execute it on map reduce

### PROCEDURE:

#### Create a sample text file

hadoop@Ubuntu:~/Documents\$ nano sample.txt

Paste the below content to sample.txt



```
GNU nano 6.2 sample.txt
1.John
2.Jane
3.Joe
4.Enna
```

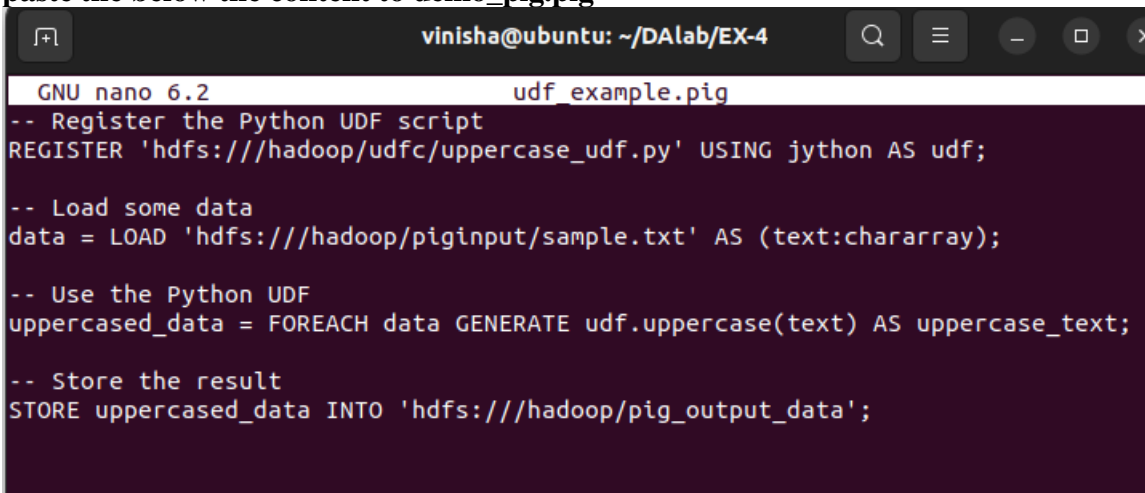
hadoop@Ubuntu:~/Documents\$ hadoop fs -put sample.txt /home/hadoop/piginput/

---

#### Create PIG File

hadoop@Ubuntu:~/Documents\$ nano demo\_pig.pig

paste the below the content to demo\_pig.pig



```
GNU nano 6.2 udf_example.pig
-- Register the Python UDF script
REGISTER 'hdfs:///hadoop/udfc/uppercase_udf.py' USING jython AS udf;

-- Load some data
data = LOAD 'hdfs:///hadoop/piginput/sample.txt' AS (text:chararray);

-- Use the Python UDF
uppercased_data = FOREACH data GENERATE udf.uppercase(text) AS uppercase_text;

-- Store the result
STORE uppercased_data INTO 'hdfs:///hadoop/pig_output_data';
```

-- Load the data from HDFS

```
data = LOAD '/home/hadoop/piginput/sample.txt' USING PigStorage(',') AS (id:int>
```

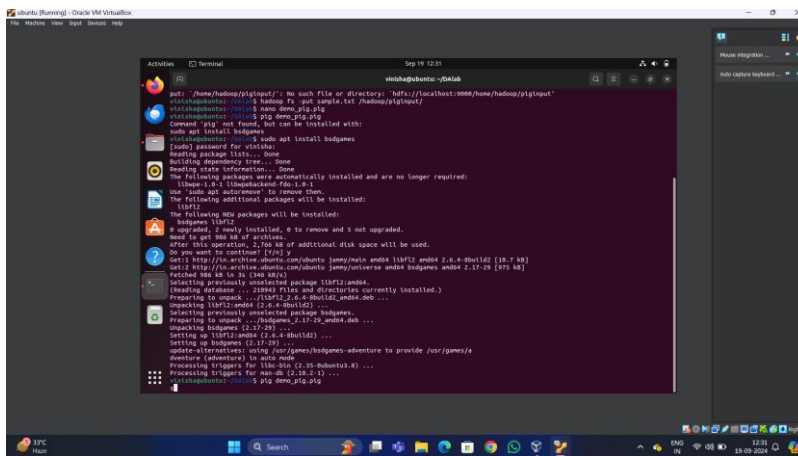
-- Dump the data to check if it was loaded correctly

DUMP data;

----- Run

**the above file**

hadoop@Ubuntu:~/Documents\$ pig demo\_pig.pig



**Create udf file an save as uppercase\_udf.py**

uppercase\_udf.py

```
def uppercase(text): return text.upper()
```

```
if __name__ == "__main__":
```

```
import sys for line in
```

```
sys.stdin:
```

```
    line = line.strip() result =
```

```
    uppercase(line)
```

```
    print(result)
```

**Create the udfs folder on hadoop**

hadoop@Ubuntu:~/Documents\$ hadoop fs -mkdir /home/hadoop/udfs

**put the upppercase\_udf.py in to the abv folder**

hadoop@Ubuntu:~/Documents\$ hdfs dfs -put uppercase\_udf.py /home/hadoop/udfs/

## To view the output

**hadoop@Ubuntu:~/Documents\$ hdfs dfs -cat /home/hadoop/pig\_output\_data/part-m00000**

```
vinisha@ubuntu:~/DALab$ nano sample.txt
vinisha@ubuntu:~/DALab$ hdfs dfs -cat /hadoop/pig_output_data/part-m-00000
1.JOHN
2.JANE
3.JOE
4.EMMA
vinisha@ubuntu:~/DALab$
```

**Result:**

Thus the program to create User Define Function in Apache Pig and execute it on map reduce has been done successfully.