## Create UDF (User Defined Functions) in Apache Pig and execute it inMapReduce/HDFS mode

### AIM:

To create UDF (User Defined Functions) in Apache Pig and execute it in MapReduce/HDFS mode.

### PROCEDURE:

1. Ensure that Apache Pig is installed and configured.



2. Create a python UDF (User Defined Functions).



3. Jython should be installed as Pig will use it to interpret the Python UDFs.

4. Create a Pig script that registers and uses the Python UDF.

```
-- Register the Python UDF script
REGISTER 'hdfs:///pig/uppercase_udf.py' USING jython AS udf;

-- Load some data
data = LOAD 'hdfs:///pig/input.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:///pig/pig_output_data';
```

5. Execute the Pig Script in MapReduce Mode using the command:
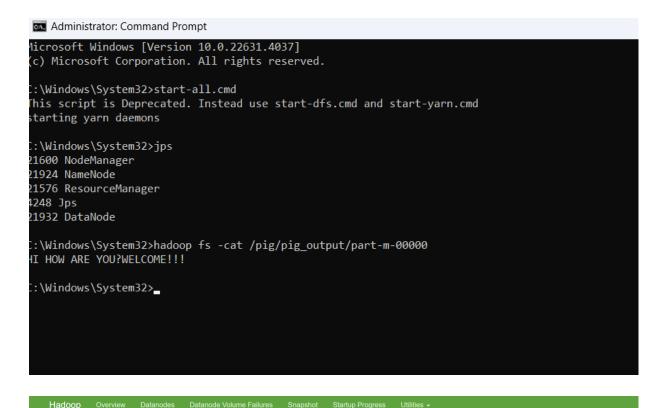
pig -x mapreduce script.pig

**OUTPUT:**

```
Administrator: Command Prompt

Microsoft Windows [Version 10.0.22631.4037]
(c) Microsoft Corporation. All rights reserved.

C:\Windows\System32>start-all.cmd
This script is Deprecated. Instead use start-dfs.cmd and start-yarn.cmd
starting yarn daemons

C:\Windows\System32>jps
21600 NodeManager
21924 NameNode
21576 ResourceManager
4248 Jps
21932 DataNode

C:\Windows\System32>hadoop fs -cat /pig/pig_output/part-m-00000
HI HOW ARE YOU?WELCOME!!!

C:\Windows\System32>_
```

Hadoop    Overview    Datanodes    Datanode Volume Failures    Snapshot    Startup Progress    Utilities ▾

# Browse Directory

/pig        Go!

Show 25 ∨ entries        Search:

| | Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name | |
|---|---|---|---|---|---|---|---|---|---|
| ☐ | -rw-r--r-- | MONASHREE | supergroup | 25 B | Sep 08 11:54 | 1 | 128 MB | input.txt | 🗑 |
| ☐ | drwxr-xr-x | MONASHREE | supergroup | 0 B | Sep 08 18:13 | 0 | 0 B | pig_output | 🗑 |
| ☐ | -rw-r--r-- | MONASHREE | supergroup | 186 B | Sep 08 18:11 | 1 | 128 MB | uppercase_udf.py | 🗑 |

Showing 1 to 3 of 3 entries        Previous 1 Next

Block ID: 1073741957

Block Pool ID: BP-468074218-172.16.11.91-1722913037250

Generation Stamp: 1133

Size: 25

Availability:

- LAPTOP-VG848917

## File contents

hi how are you?welcome!!!

Close

Block information --   Block 0 ▾

Block ID: 1073742014

Block Pool ID: BP-468074218-172.16.11.91-1722913037250

Generation Stamp: 1190

Size: 26

Availability:

- LAPTOP-VG848917

**File contents**

HI HOW ARE YOU?WELCOME!!!

**RESULT:**

Thus, to create a UDF in Apache Pig and execute in MapReduce mode has been executed successfully