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

## AIM:

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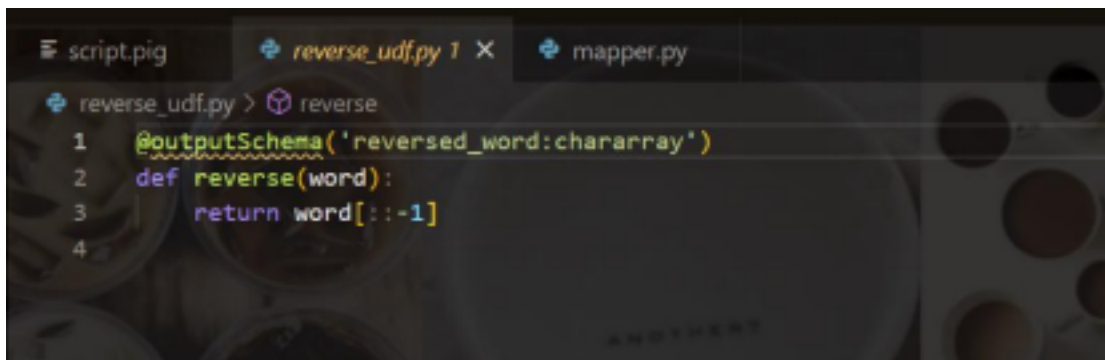
## PROCEDURE:

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.

5. Execute the Pig Script in MapReduce Mode using the command:

pig -x mapreduce script.pig

## OUTPUT:





```
C:\>hadoop fs -mkdir /pig

C:\>hadoop fs -put C:/Users/mercy/OneDrive/Documents/DataAnalytics/Pig/input.txt /pig

C:\>hadoop fs -cat /pig/output/part-m-00000
olleH
avaJ
ycreM
nohtyP
gnaM
etalocohC
eeffoC
```

## RESULT:

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