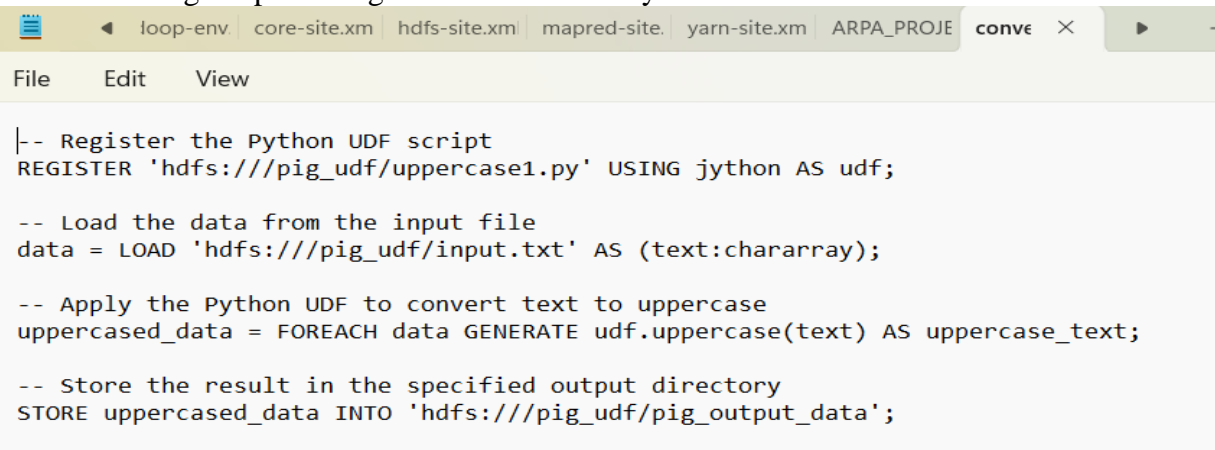## Ex:5   Create UDF (User Defined Functions) in Apache Pig
## and execute it inMapReduce/HDFS mode

## 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_udf/uppercase1.py' USING jython AS udf;

-- Load the data from the input file
data = LOAD 'hdfs:///pig_udf/input.txt' AS (text:chararray);

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

-- Store the result in the specified output directory
STORE uppercased_data INTO 'hdfs:///pig_udf/pig_output_data';
```

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

pig -x mapreduce script.pig

## OUTPUT:
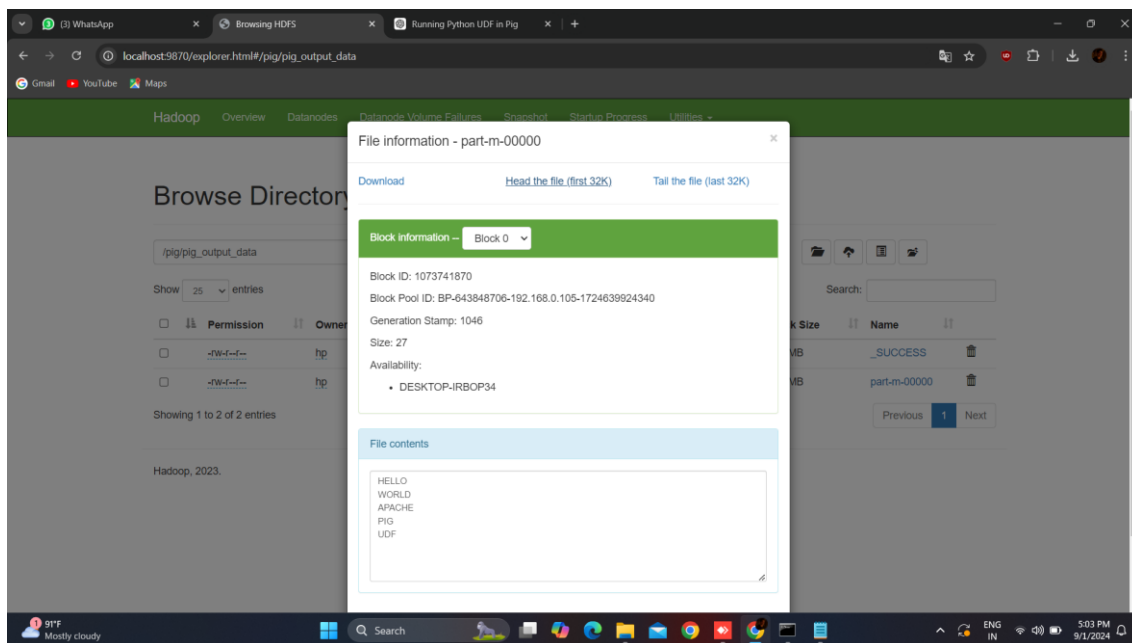
**RESULT:**

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