**Bharathwin MA**
**205229105**
**Big Data Management and Analytics**
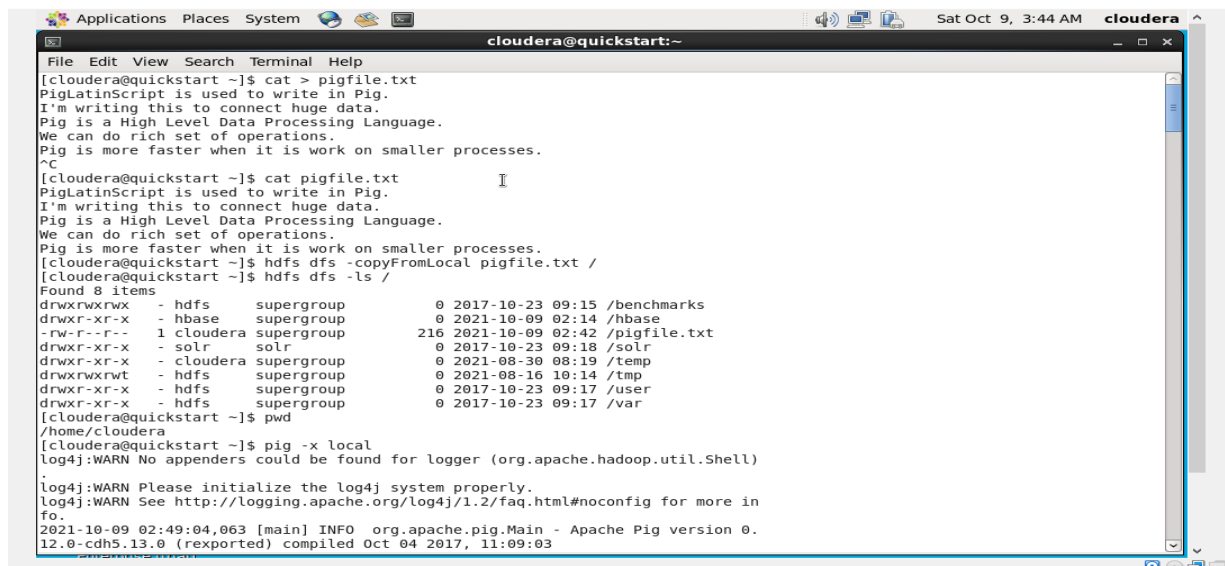

<div align="center">

**Exercise 07: Word Count using Pig Grouping**

</div>

Here, we will be running Apache Pig Sample scripts using grunts. It is to just  see the power of Apache Pig.
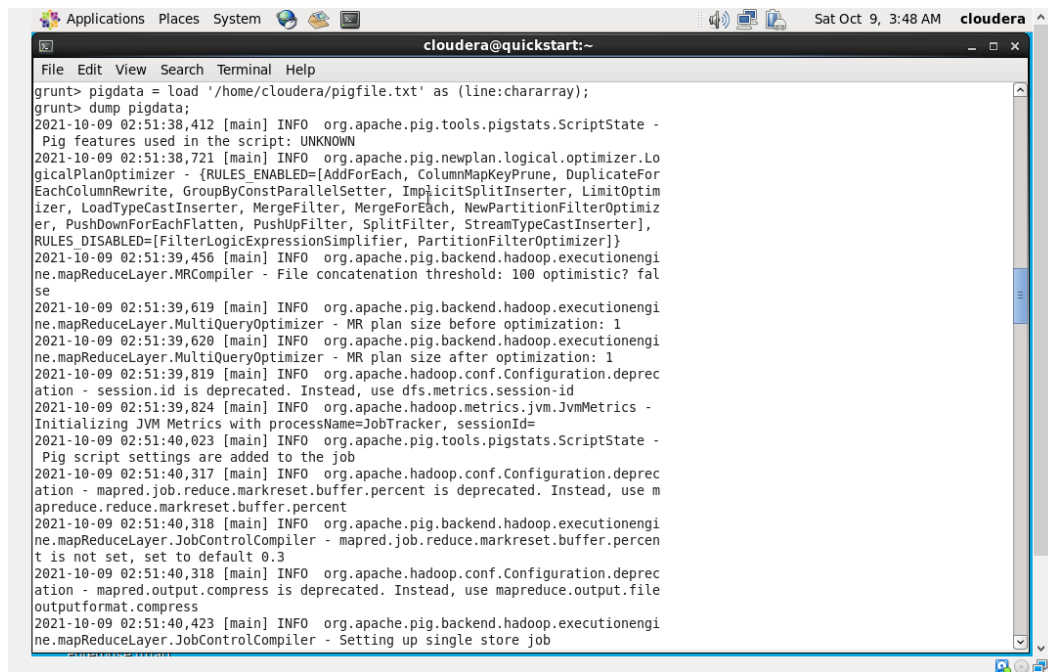


**Step 1A: Start Grunt shell.**

Open terminal and type *pig*



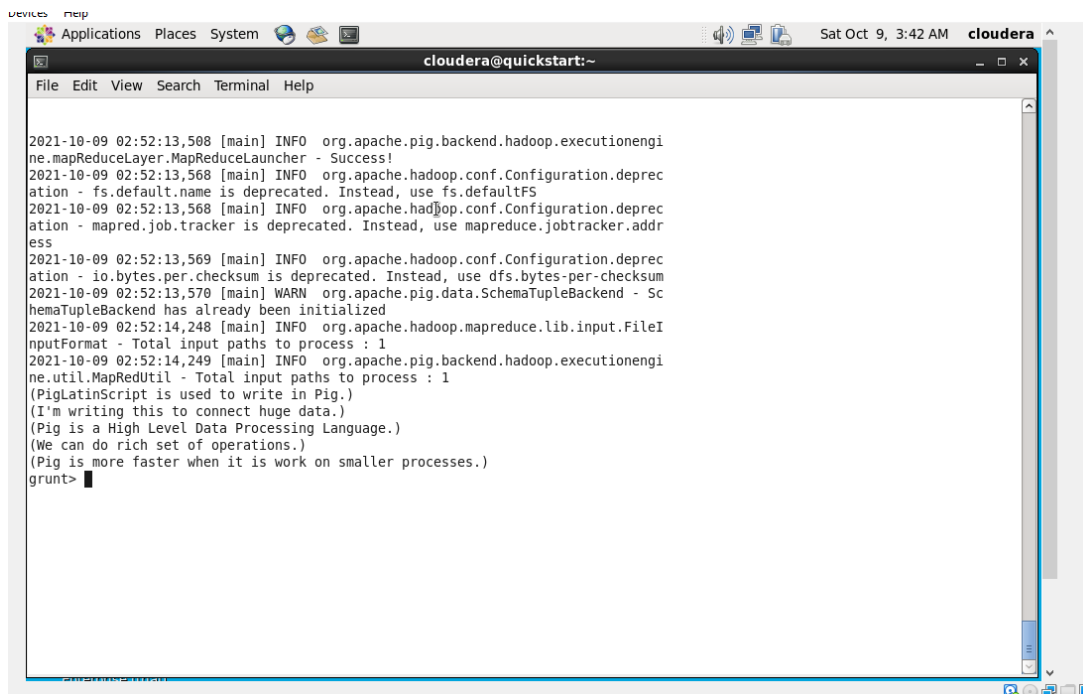**Step 1B: Create a file at /user/cloudera/pigfile.txt**

## Step 2 : Load the file stored in hdfs and each line have to store in 'line' (Space separated file)

## Step 3: flatten the words in each line and save separated words into a variable

*grunt>wordsinline = FOREACH input1 GENERATE flatten(TOKENIZE(line, ' ')) as word;*
*grunt>DUMP wordsinline;*

## Step 4: Group the similar words and save into a variable

*grunt>groupwords = _____ wordsinline by word;*
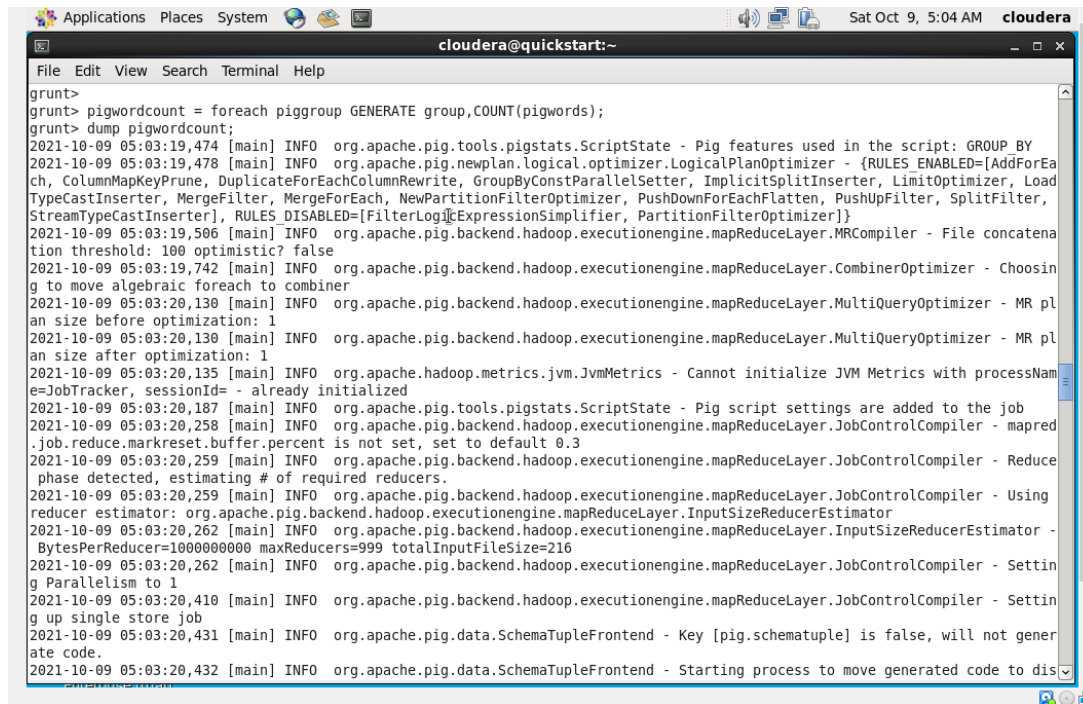*grunt>dump groupwords;*
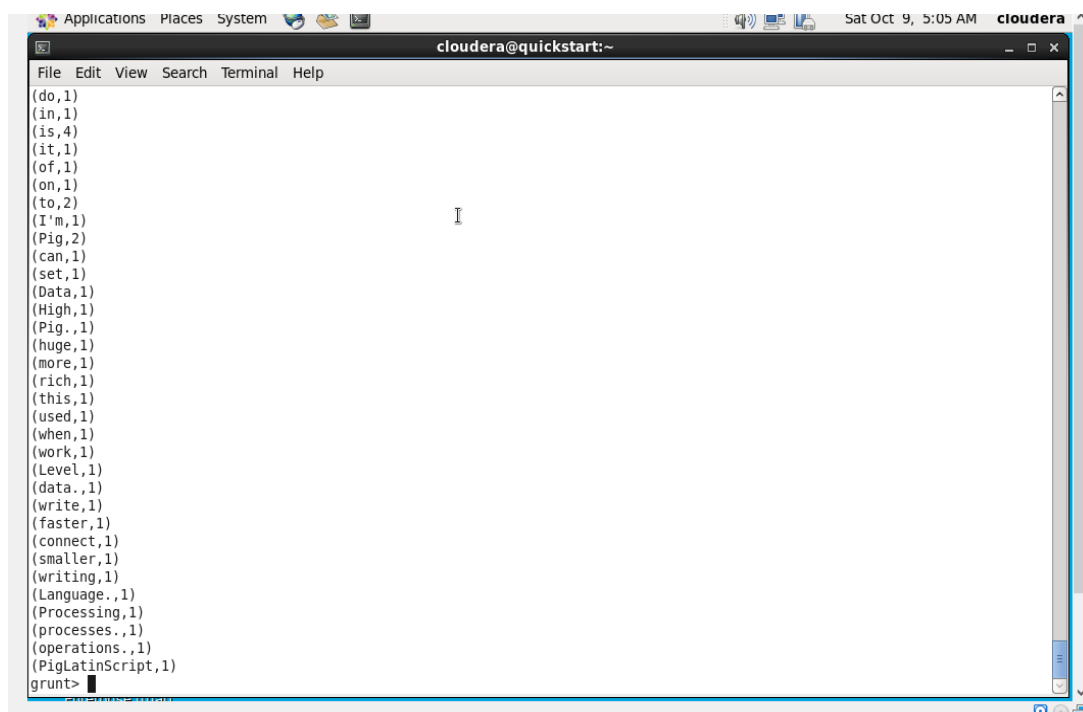*grunt>describe groupwords;*

## Step 5: Count Words in the group.

*grunt>countwords = foreach _____;*
*grunt>DUMP countwords;*