

File Edit View VM Tabs Help | || | | | | | | |

Ubuntu 64-bit X

Activities Terminal ▾

Nov 19 18:26

Speaker Power

```
hadoop@ubuntu:~$ hdfs dfs -put genderClassification.csv /TestData/genderClassification.csv
hadoop@ubuntu:~$ hdfs dfs -ls /TestData/genderClassification.csv
-rw-r--r-- 1 hadoop supergroup 2418 2021-11-19 18:20 /TestData/genderClassification.csv
hadoop@ubuntu:~$ hdfs dfs -cat /TestData/genderClassification.csv
Favorite Color,Favorite Music Genre,Favorite Beverage,Favorite Soft Drink,Gender
Cool,Rock,Vodka,7UP/Sprite,F
Neutral,Hip hop,Vodka,Coca Cola/Pepsi,F
Warm,Rock,Wine,Coca Cola/Pepsi,F
Warm,Folk/Traditional,Whiskey,Fanta,F
Cool,Rock,Vodka,Coca Cola/Pepsi,F
Warm,Jazz/Blues,Doesn't drink,Fanta,F
Cool,Pop,Beer,Coca Cola/Pepsi,F
Warm,Pop,Whiskey,Fanta,F
Warm,Rock,Other,7UP/Sprite,F
Neutral,Pop,Wine,Coca Cola/Pepsi,F
Cool,Pop,Other,7UP/Sprite,F
Warm,Pop,Other,7UP/Sprite,F
Warm,Pop,Wine,7UP/Sprite,F
Warm,Electronic,Wine,Coca Cola/Pepsi,F
Cool,Rock,Beer,Coca Cola/Pepsi,F
Warm,Jazz/Blues,Wine,Coca Cola/Pepsi,F
Cool,Pop,Wine,7UP/Sprite,F
Cool,Rock,Other,Coca Cola/Pepsi,F
Cool,Rock,Other,Coca Cola/Pepsi,F
Cool,Pop,Doesn't drink,7UP/Sprite,F
Cool,Pop,Beer,Fanta,F
Warm,Jazz/Blues,Whiskey,Fanta,F
Cool,Rock,Vodka,Coca Cola/Pepsi,F
Warm,Pop,Other,Coca Cola/Pepsi,F
```

*Untitled - Notepad

File Edit Format View Help
Student Name = Bibek Shah Shankhar

1. We copied genderClassification.csv from local system to Hdfs.
- hdfs dfs -put genderClassification.csv /TestData/genderClassification.csv
2. We listed genderClassification.csv to confirm whether the file exists in hdfs or not.
- hdfs dfs -ls genderClassification.csv /TestData/genderClassification.csv
3. Then We print the 'genderClassification' file to know the file
- hdfs dfs -cat genderClassification.csv /TestData/genderClassification.csv

Ln 15, Col 29

100%

Windows (CRLF)

UTF-8

To direct input to this VM, move the mouse pointer inside or press Ctrl+G.



8:11 AM 20-Nov-21 1

File Edit View VM Tabs Help | || | | | | | | | |

Ubuntu 64-bit X

Activities Terminal ▾

Nov 19 18:55

hadoop@ubuntu:~\$ pig
21/11/19 18:54:40 INFO pig.ExecTypeProvider: Trying ExecType : LOCAL
21/11/19 18:54:40 INFO pig.ExecTypeProvider: Trying ExecType : MAPREDUCE
21/11/19 18:54:40 INFO pig.ExecTypeProvider: Picked MAPREDUCE as the ExecType
2021-11-19 18:54:40,259 [main] INFO org.apache.pig.Main - Apache Pig version 0.15.0 (r1682971) co
mpiled Jun 01 2015, 11:44:35
2021-11-19 18:54:40,259 [main] INFO org.apache.pig.Main - Logging error messages to: /home/hadoop
/pig_1637376880257.log
2021-11-19 18:54:40,302 [main] INFO org.apache.pig.impl.util.Utils - Default bootup file /home/ha
dooop/.pigbootup not found
2021-11-19 18:54:40,674 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job
.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-11-19 18:54:40,674 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default
.name is deprecated. Instead, use fs.defaultFS
2021-11-19 18:54:40,675 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngin
e - Connecting to hadoop file system at: hdfs://localhost:8020
2021-11-19 18:54:41,407 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default
.name is deprecated. Instead, use fs.defaultFS
grunt> genderClassification = LOAD '/TestData/genderClassification.csv'
-> USING PigStorage(',')
-> AS (fav_color:chararray, fav_music_genre:chararray,
-> fav_beverage:chararray, fav_soft_drink:chararray,
-> gender:chararray);
2021-11-19 18:54:52,706 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default
.name is deprecated. Instead, use fs.defaultFS
grunt>

*Untitled - Notepad

File Edit Format View Help

Student Name = Bibek Shah Shankhar

Question 1. Execute Load command with an Example.

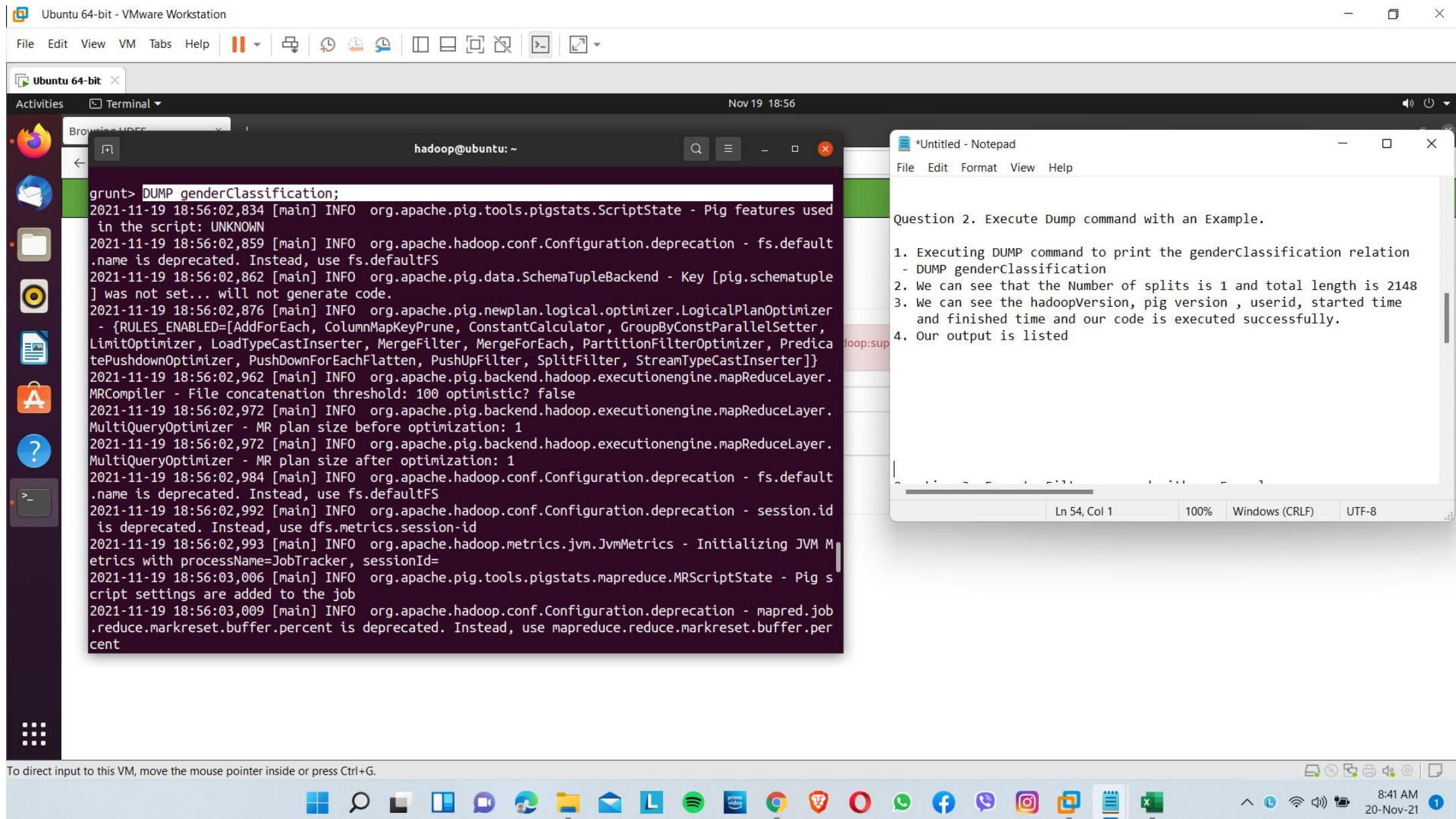
1. Starting pig grunt shell
2. Loading the 'genderClassification.csv' into relation named genderClassification .
 - genderClassification = LOAD '/TestData/genderClassification.csv'
USING PigStorage(',')
AS (fav_color:chararray, fav_music_genre:chararray,
fav_beverage:chararray, fav_soft_drink:chararray,
gender:chararray);

Ln 33, Col 1 100% Windows (CRLF) UTF-8

To direct input to this VM, move the mouse pointer inside or press Ctrl+G.

8:40 AM 20-Nov-21





File Edit View VM Tabs Help | || | | | | | | | |

Ubuntu 64-bit X

Activities Terminal ▾

Nov 19 18:56

```
hadoop@ubuntu:~$ hadoop fs -ls /user/hadoop/genderClassification/part-r-00000
2021-11-19 18:56:04,866 [Thread-18] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-11-19 18:56:04,869 [Thread-18] INFO org.apache.hadoop.mapred.LocalJobRunner - OutputCommitter is org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.PigOutputCommitter
2021-11-19 18:56:04,897 [Thread-18] INFO org.apache.hadoop.mapred.LocalJobRunner - Waiting for map tasks
2021-11-19 18:56:04,897 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.LocalJobRunner - Starting task: attempt_local59974511_0001_m_000000_0
2021-11-19 18:56:04,947 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.Task - Using ResourceCalculatorProcessTree : [ ]
2021-11-19 18:56:04,950 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.MapTask - Processing split: Number of splits :1
Total Length = 2418
Input split[0]:
Length = 2418
ClassName: org.apache.hadoop.mapreduce.lib.input.FileSplit
Locations:
-----
2021-11-19 18:56:04,964 [LocalJobRunner Map Task Executor #0] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.PigRecordReader - Current split being processed hdfs://localhost:8020/TestData/genderClassification.csv:0+2418
2021-11-19 18:56:05,001 [LocalJobRunner Map Task Executor #0] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-11-19 18:56:05,019 [LocalJobRunner Map Task Executor #0] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.PigMapOnly$Map - Aliases being processed per job phase (AliasName[line,offset]): M: genderClassification[1,23],genderClassification[-1,-1] C: R:
2021-11-19 18:56:05,048 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.Local
```

*Untitled - Notepad

File Edit Format View Help

Question 2. Execute DUMP command with an Example.

1. Executing DUMP command to print the genderClassification relation
- DUMP genderClassification
2. We can see that the Number of splits is 1 and total length is 2148
3. We can see the hadoopVersion, pig version , userid, started time and finished time and our code is executed successfully.
4. Our output is listed

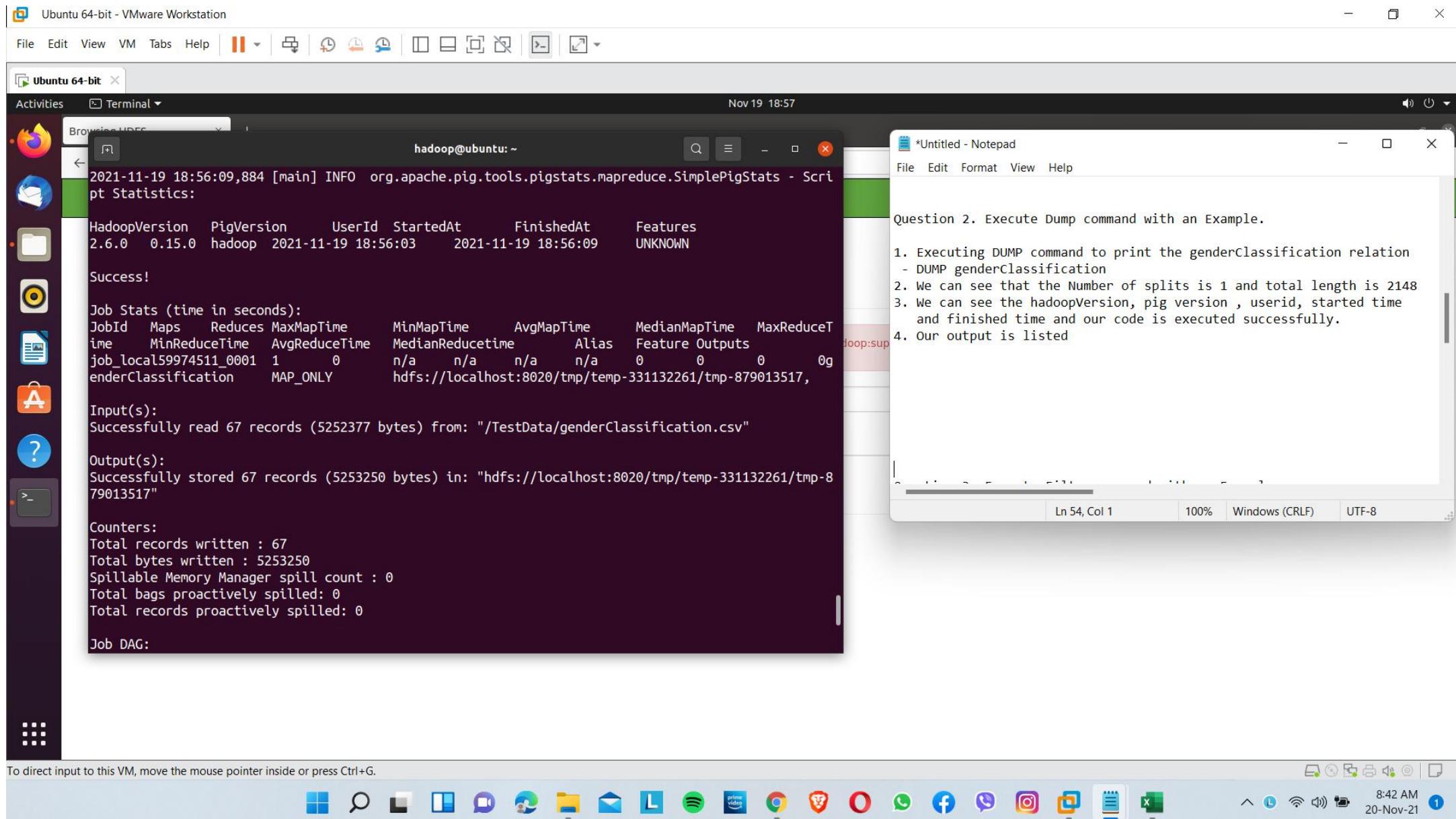
Ln 54, Col 1

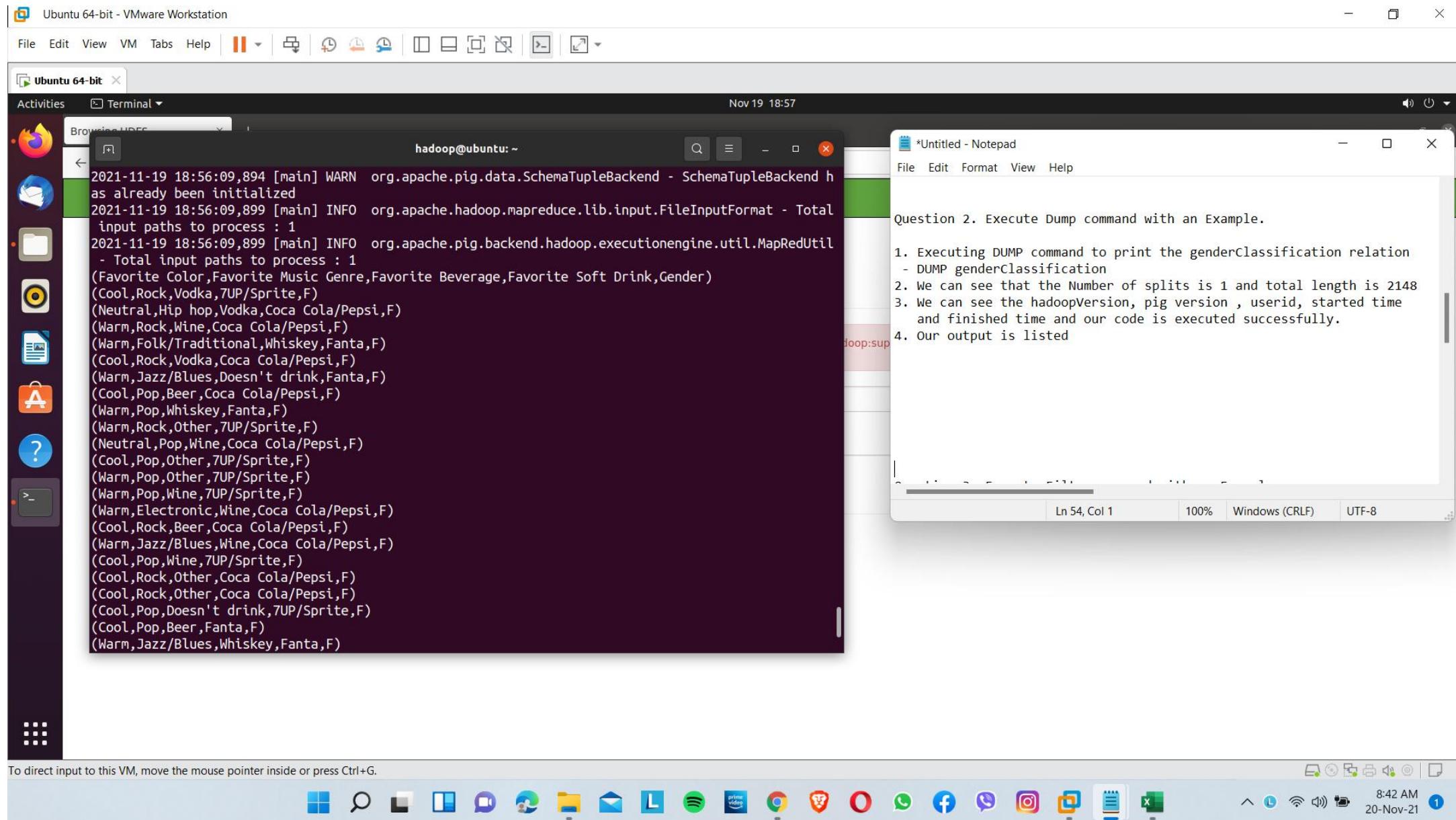
100% Windows (CRLF) UTF-8

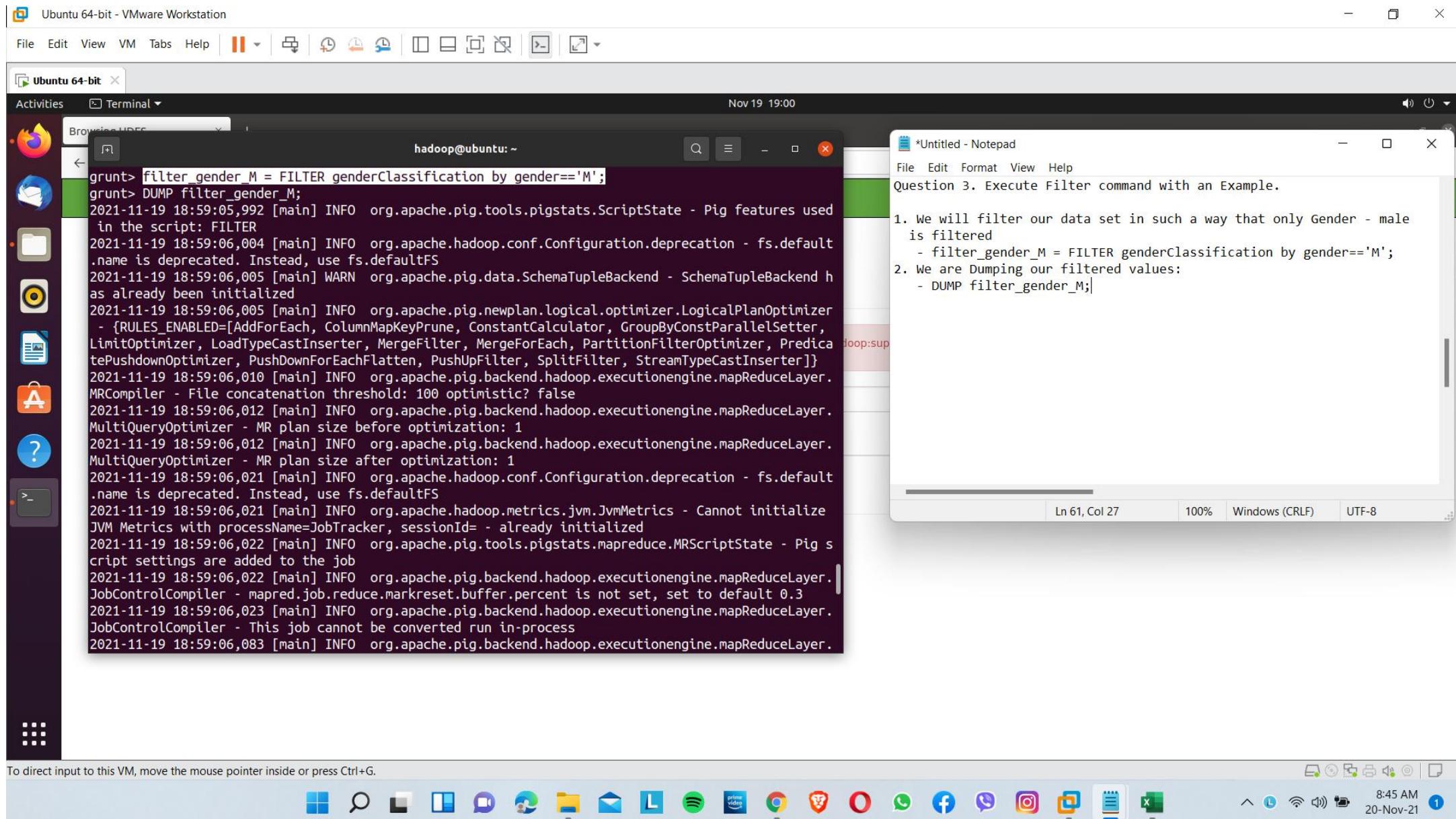
To direct input to this VM, move the mouse pointer inside or press Ctrl+G.

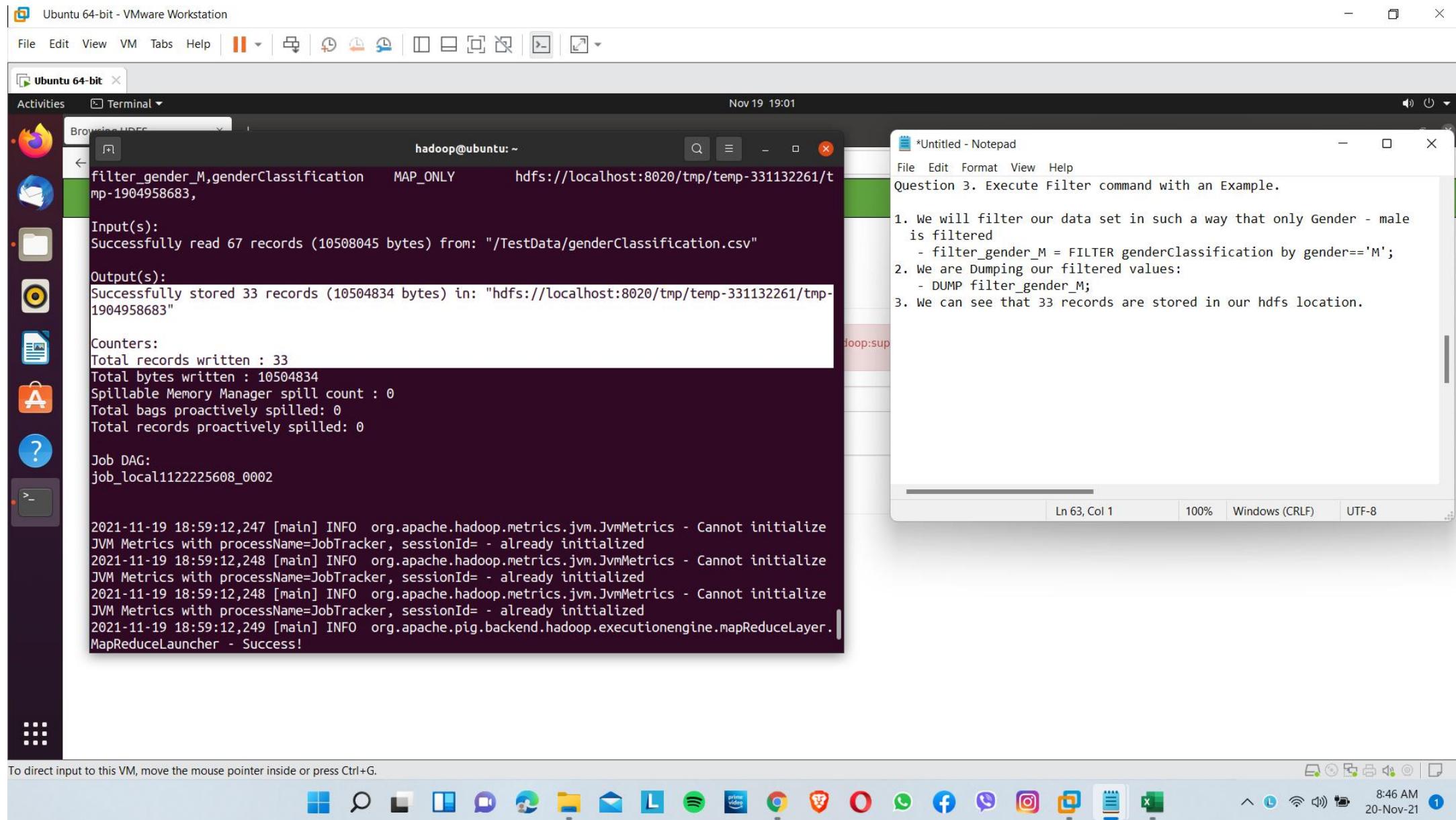


8:41 AM 20-Nov-21 1









File Edit View VM Tabs Help | || | | | | | | | |

Ubuntu 64-bit X

Activities Terminal ▾

Nov 19 19:02

Browsing HDFS

```
hadoop@ubuntu:~
```

```
- Total input paths to process : 1
(Warm,R&B and soul,Wine,Other,M)
(Neutral,Hip hop,Beer,7UP/Sprite,M)
(Warm,Electronic,Other,Coca Cola/Pepsi,M)
(Neutral,Rock,Doesn't drink,Coca Cola/Pepsi,M)
(Cool,Pop,Other,Fanta,M)
(Cool,Pop,Whiskey,Fanta,M)
(Warm,Rock,Vodka,7UP/Sprite,M)
(Cool,Rock,Vodka,Coca Cola/Pepsi,M)
(Neutral,Pop,Doesn't drink,7UP/Sprite,M)
(Warm,R&B and soul,Doesn't drink,Coca Cola/Pepsi,M)
(Cool,Rock,Wine,7UP/Sprite,M)
(Cool,Folk/Traditional,Beer,Other,M)
(Cool,Hip hop,Beer,Coca Cola/Pepsi,M)
(Cool,Hip hop,Wine,Coca Cola/Pepsi,M)
(Cool,R&B and soul,Whiskey,7UP/Sprite,M)
(Cool,Rock,Doesn't drink,Other,M)
(Warm,Hip hop,Beer,Coca Cola/Pepsi,M)
(Cool,R&B and soul,Doesn't drink,Coca Cola/Pepsi,M)
(Cool,Rock,Doesn't drink,Coca Cola/Pepsi,M)
(Cool,Hip hop,Doesn't drink,Other,M)
(Warm,Rock,Beer,Fanta,M)
(Cool,Electronic,Doesn't drink,Fanta,M)
(Cool,Electronic,Other,Fanta,M)
(Warm,Folk/Traditional,Other,Fanta,M)
(Warm,Electronic,Vodka,Fanta,M)
(Warm,Jazz/Blues,Vodka,Coca Cola/Pepsi,M)
(Cool,Pop,Whiskey,Other,M)
(Cool,Electronic,Whiskey,Coca Cola/Pepsi,M)
```

*Untitled - Notepad

File Edit Format View Help

Question 3. Execute Filter command with an Example.

1. We will filter our data set in such a way that only Gender - male is filtered
 - filter_gender_M = FILTER genderClassification by gender=='M';
2. We are Dumping our filtered values:
 - DUMP filter_gender_M;
3. We can see that 33 records are stored in our hdfs location.
4. Finally, The gender with male or M class are are filtered out

Ln 63, Col 65

100% Windows (CRLF) UTF-8

To direct input to this VM, move the mouse pointer inside or press Ctrl+G.

8:47 AM 20-Nov-21



File Edit View VM Tabs Help | || | | | | | | | |

Ubuntu 64-bit X

Activities Terminal ▾

Nov 19 19:03

```
hadoop@ubuntu:~$ grunt> filter_gender_F = FILTER genderClassification by gender=='F';
grunt> DUMP filter_gender_F;
2021-11-19 19:02:28,088 [main] INFO  org.apache.pig.tools.pigstats.ScriptState - Pig features used
in the script: FILTER
2021-11-19 19:02:28,101 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - fs.default
.name is deprecated. Instead, use fs.defaultFS
2021-11-19 19:02:28,101 [main] WARN  org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend h
as already been initialized
2021-11-19 19:02:28,101 [main] INFO  org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer
- [RULES_ENABLED=[AddForEach, ColumnMapKeyPrune, ConstantCalculator, GroupByConstParallelSetter,
LimitOptimizer, LoadTypeCastInserter, MergeFilter, MergeForEach, PartitionFilterOptimizer, Predica
tePushdownOptimizer, PushDownForEachFlatten, PushUpFilter, SplitFilter, StreamTypeCastInserter]]
2021-11-19 19:02:28,112 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.
MRCompiler - File concatenation threshold: 100 optimistic? false
2021-11-19 19:02:28,113 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.
MultiQueryOptimizer - MR plan size before optimization: 1
2021-11-19 19:02:28,113 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.
MultiQueryOptimizer - MR plan size after optimization: 1
2021-11-19 19:02:28,122 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - fs.default
.name is deprecated. Instead, use fs.defaultFS
2021-11-19 19:02:28,123 [main] INFO  org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize
JVM Metrics with processName=JobTracker, sessionId= - already initialized
2021-11-19 19:02:28,124 [main] INFO  org.apache.pig.tools.pigstats.mapreduce.MRScriptState - Pig s
cript settings are added to the job
2021-11-19 19:02:28,125 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.
JobControlCompiler - mapred.job.reduce.markreset.buffer.percent is not set, set to default 0.3
2021-11-19 19:02:28,125 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.
JobControlCompiler - This job cannot be converted run in-process
2021-11-19 19:02:28,357 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.
```

*Untitled - Notepad

File Edit Format View Help

Question 3. Execute Filter command with an Example.

1. We have filtered our data set in such a way that only Gender - male is filtered
 - filter_gender_M = FILTER genderClassification by gender=='M';
2. We are Dumping our filtered values:
 - DUMP filter_gender_M;
3. We can see that 33 records are stored in our hdfs location.
4. Finally, The gender with male or M class are are filtered out
5. We have filtered our data set in such a way that only Gender - Female is filtered
 - filter_gender_F = FILTER genderClassification by gender=='F';

Ln 67, Col 65

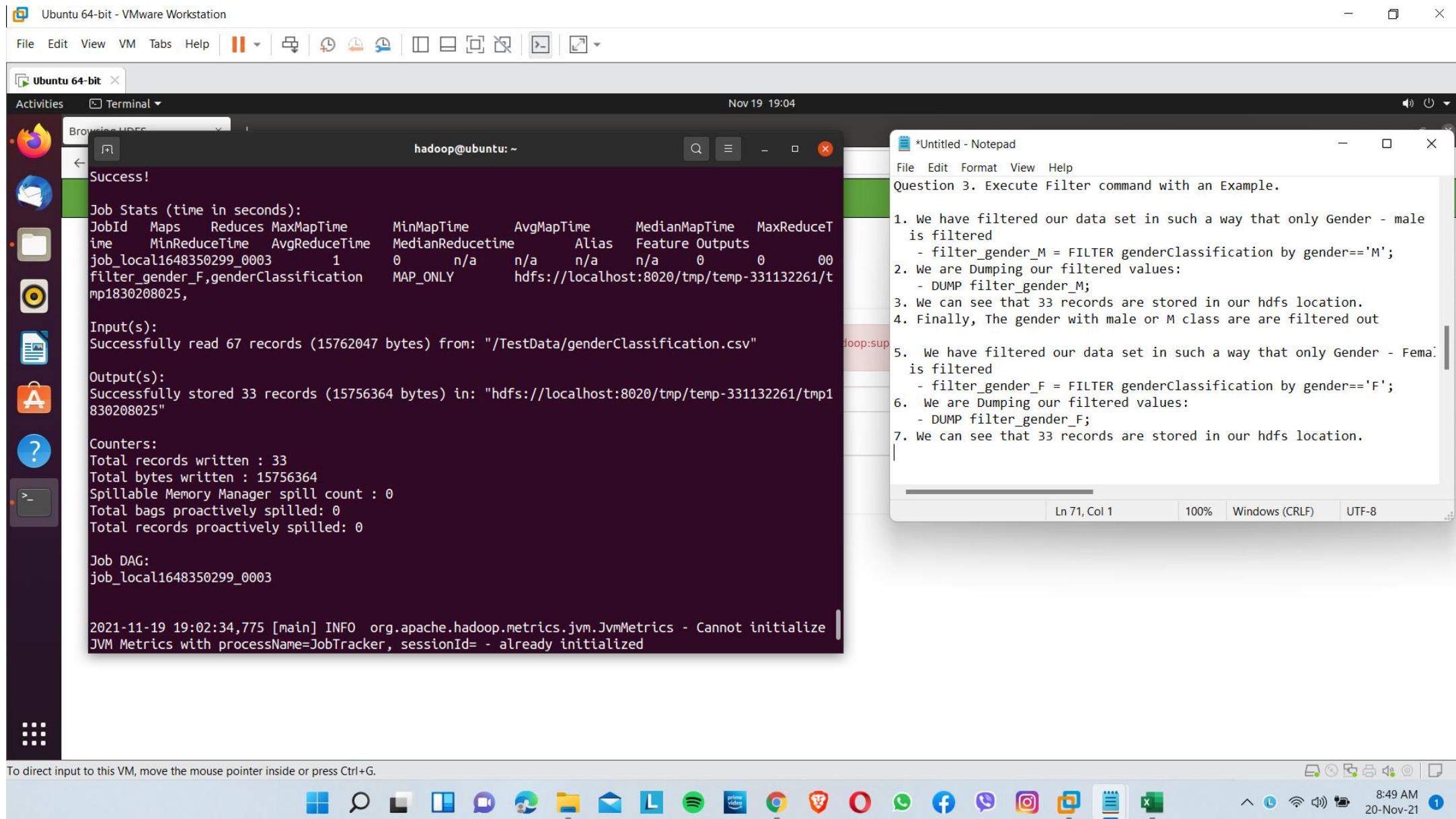
100% Windows (CRLF)

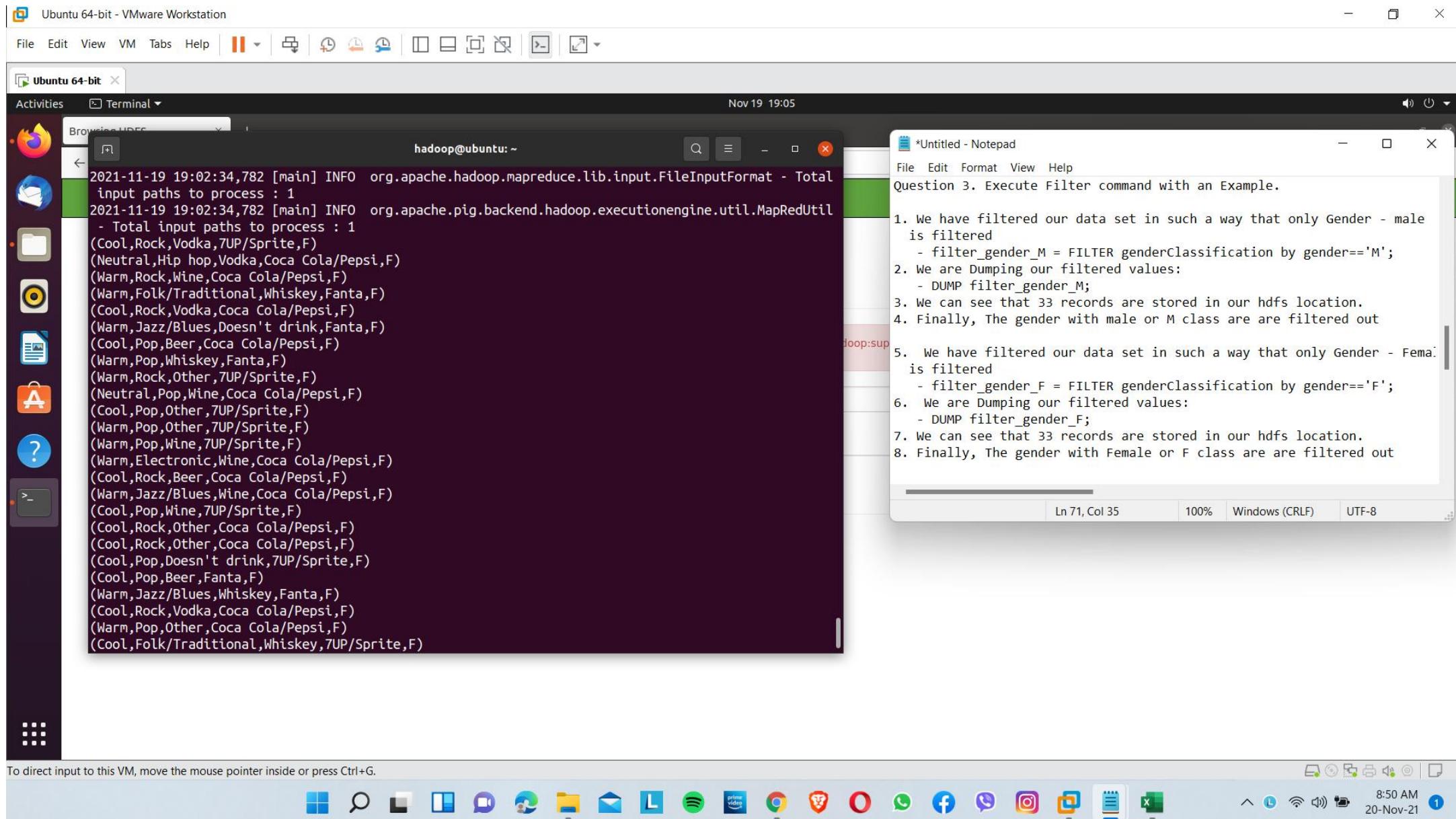
UTF-8

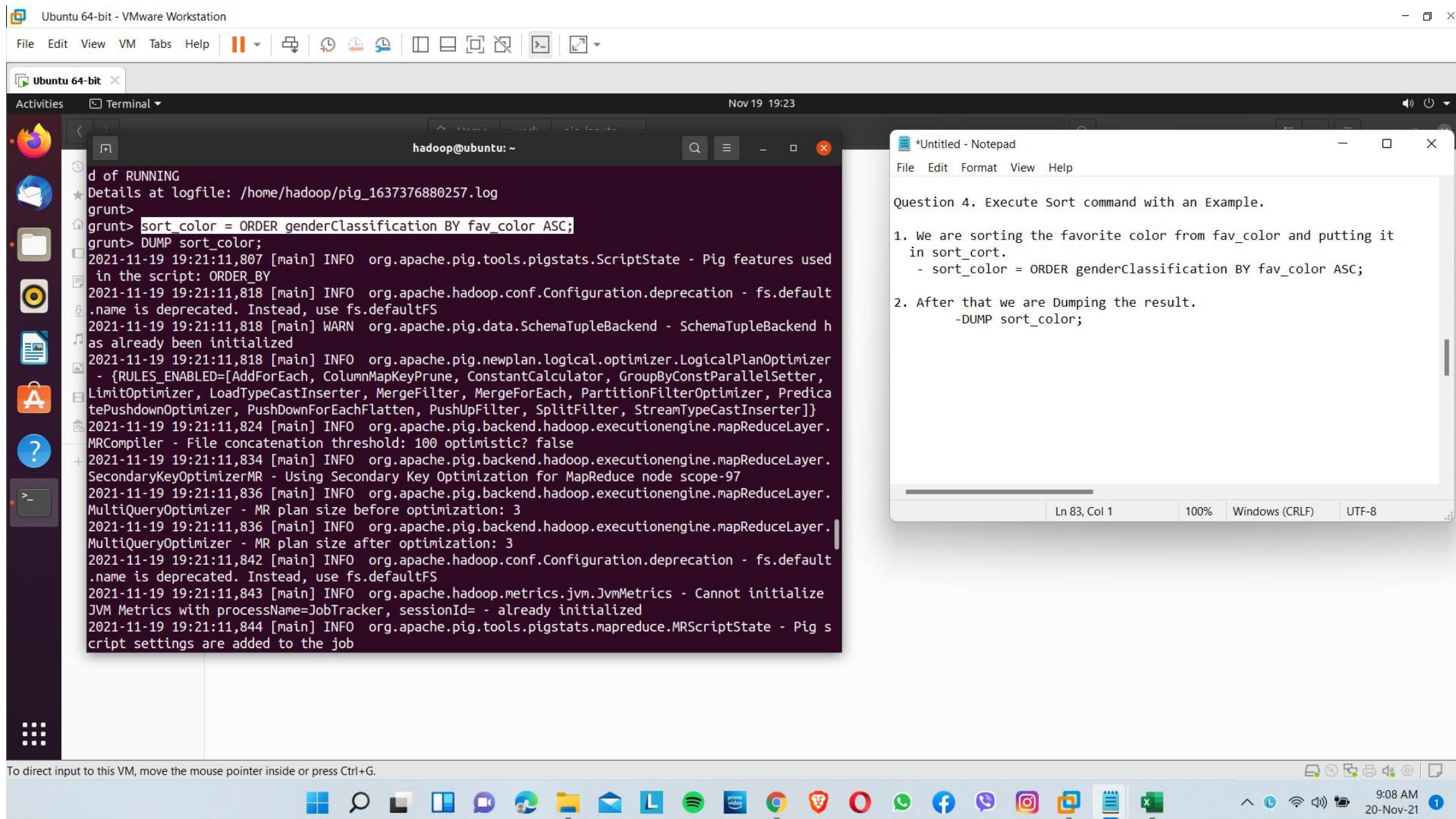
To direct input to this VM, move the mouse pointer inside or press Ctrl+G.

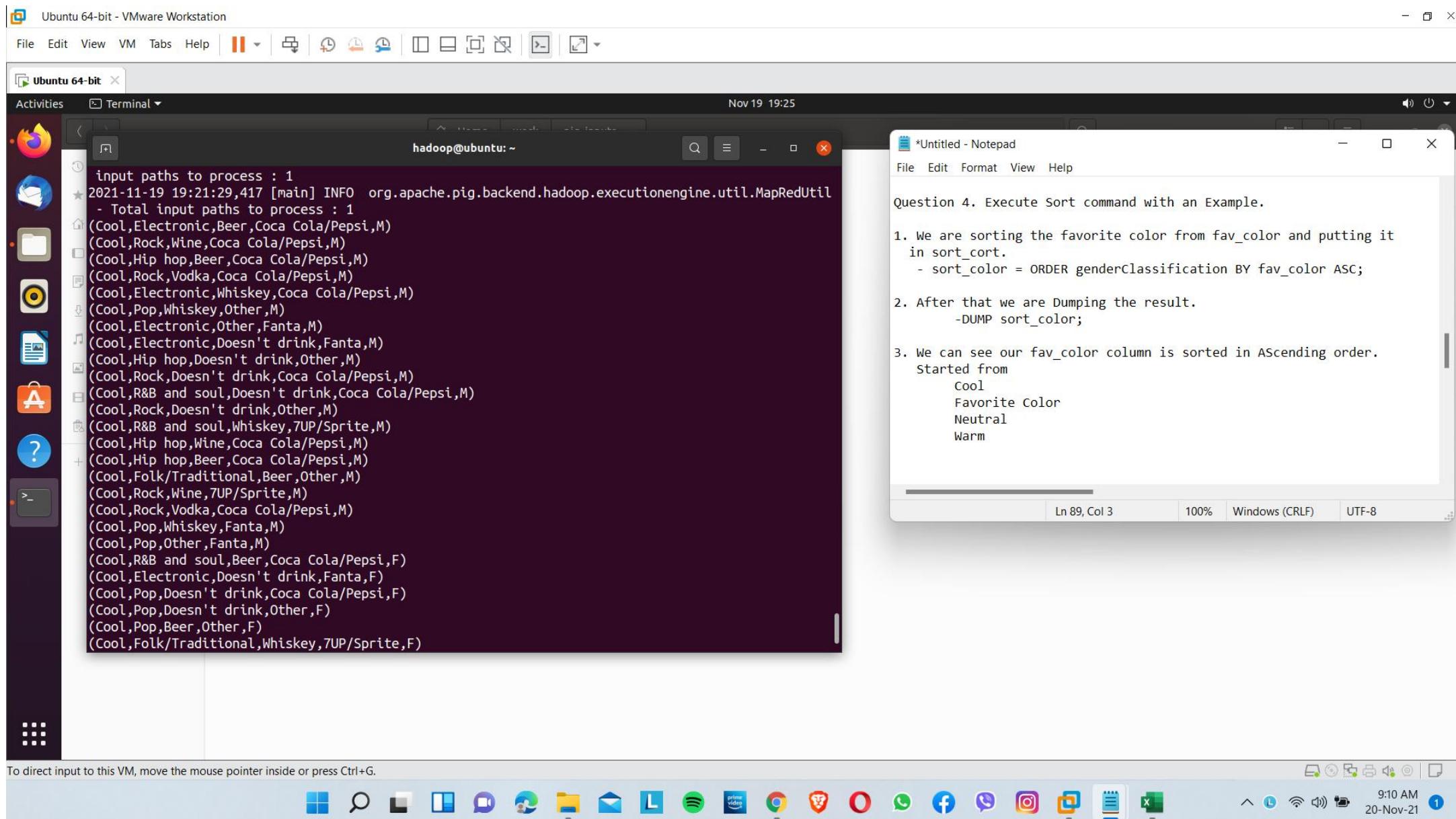


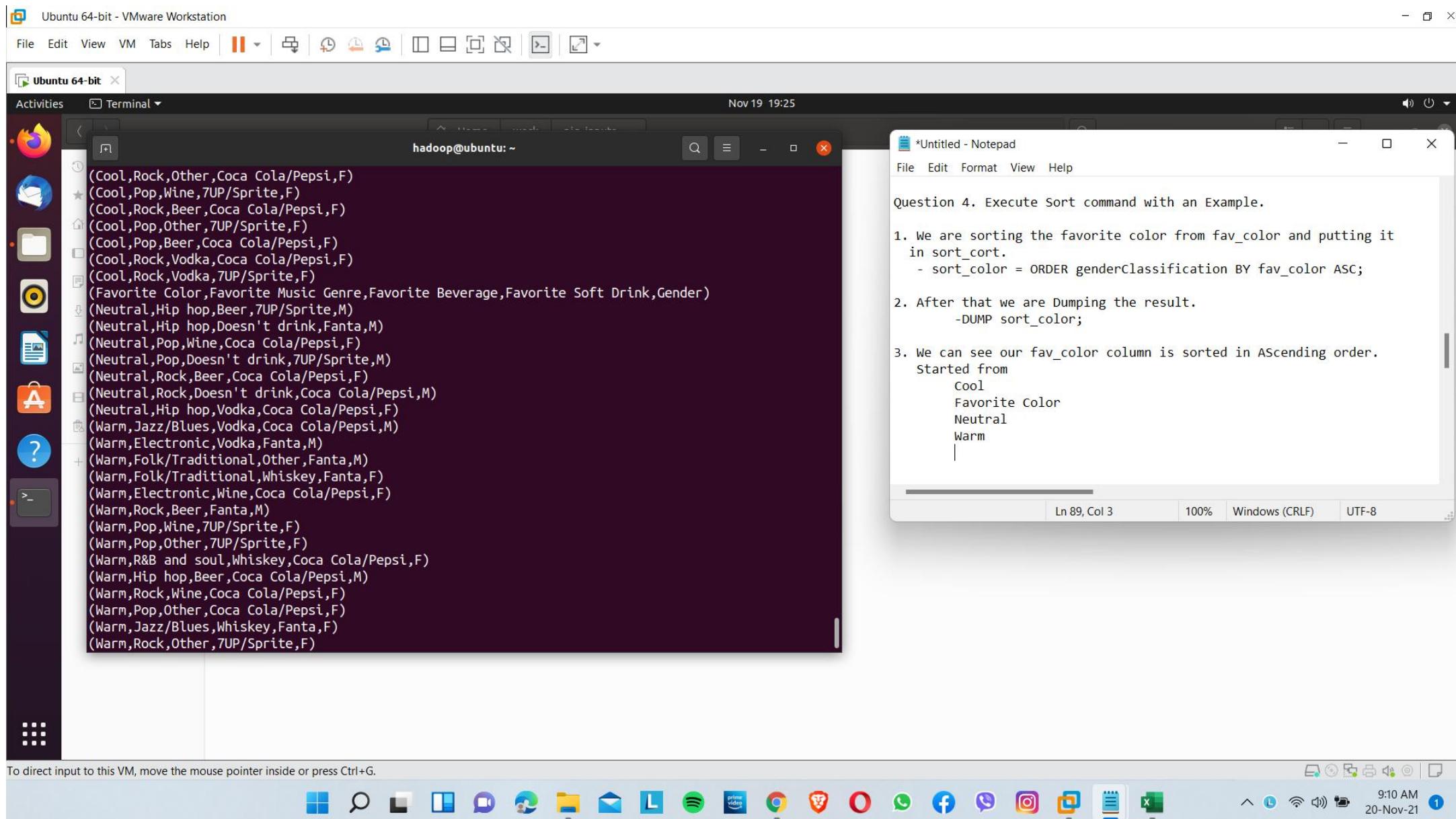
8:48 AM 20-Nov-21 1

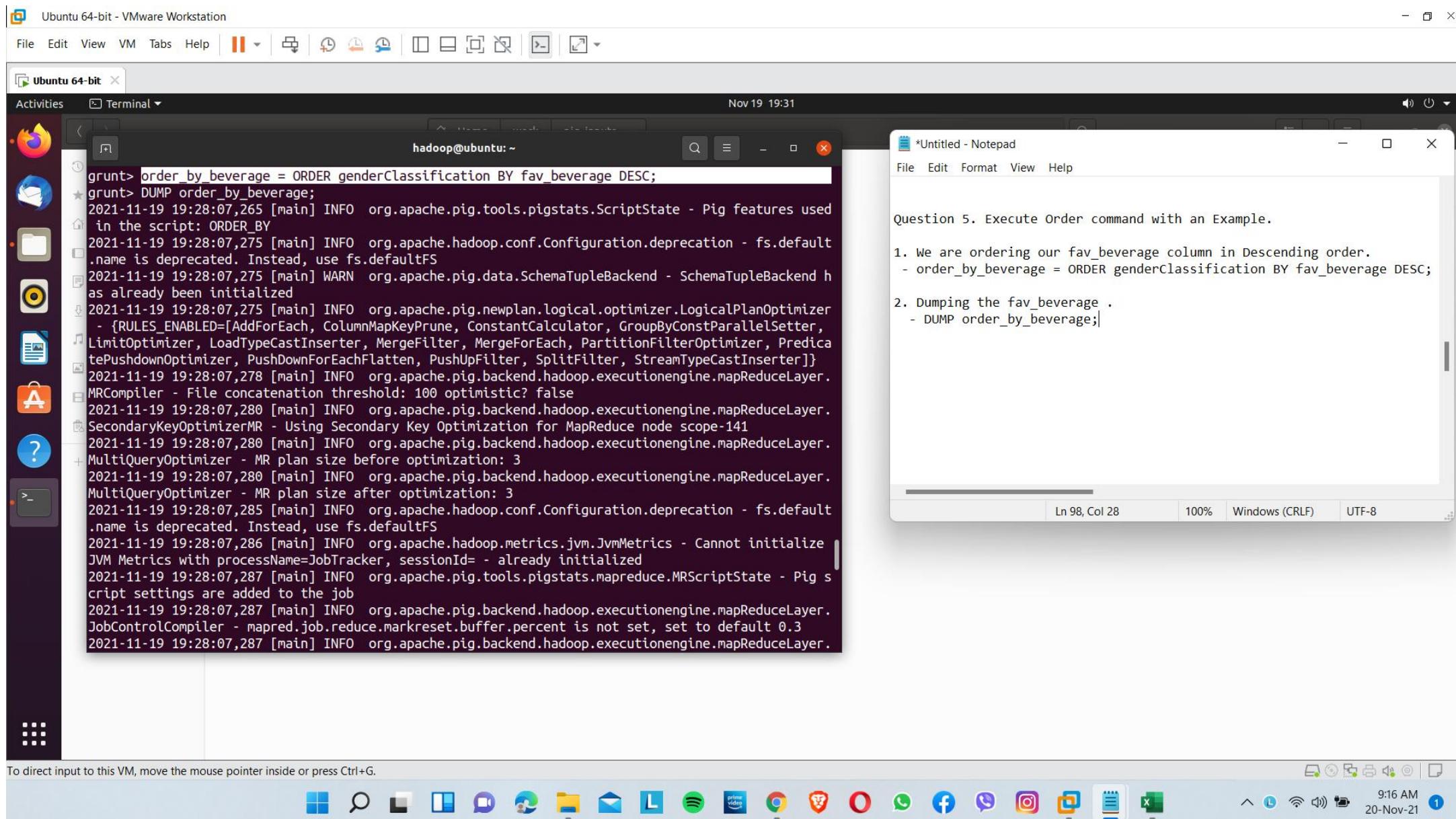


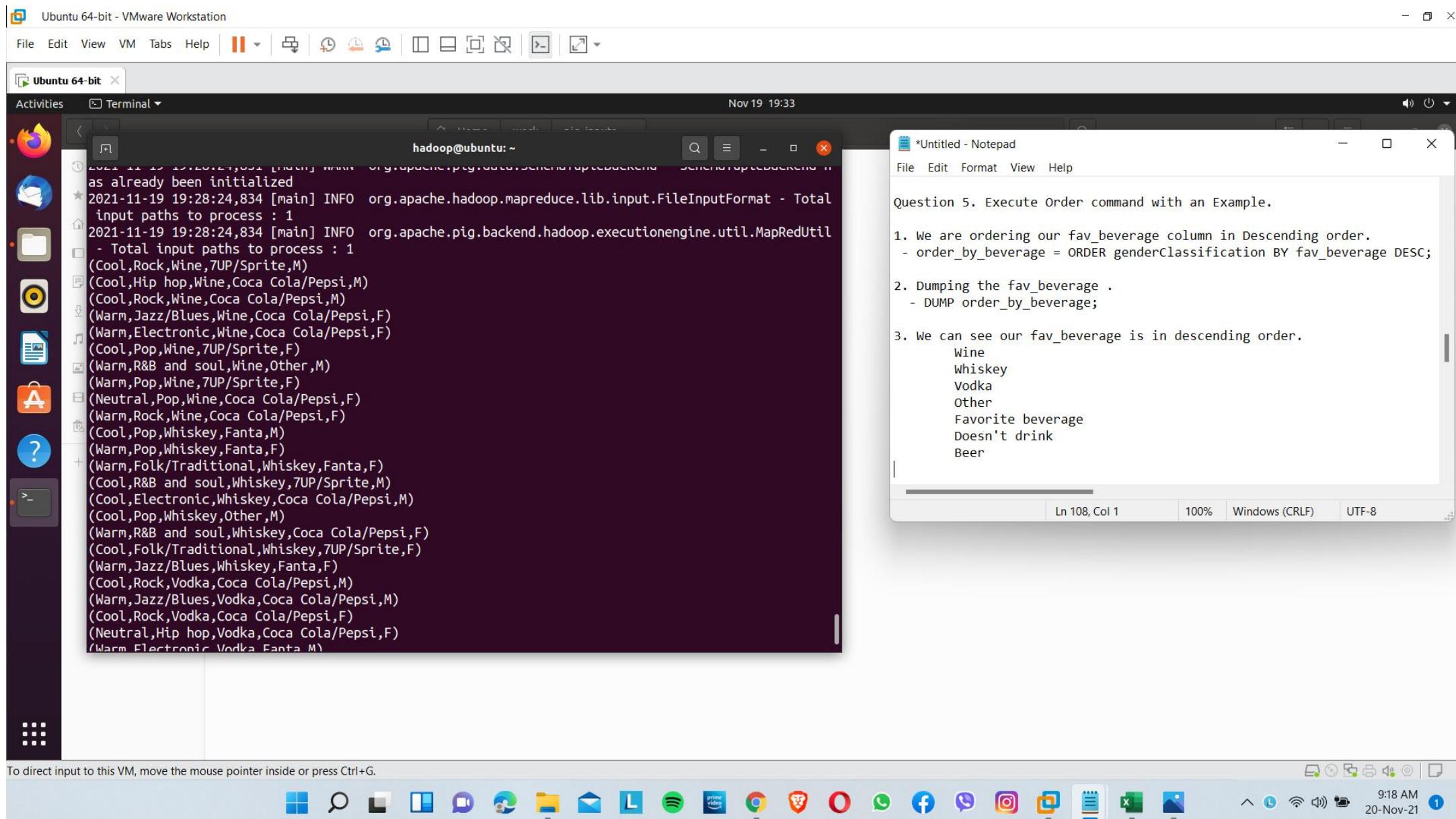


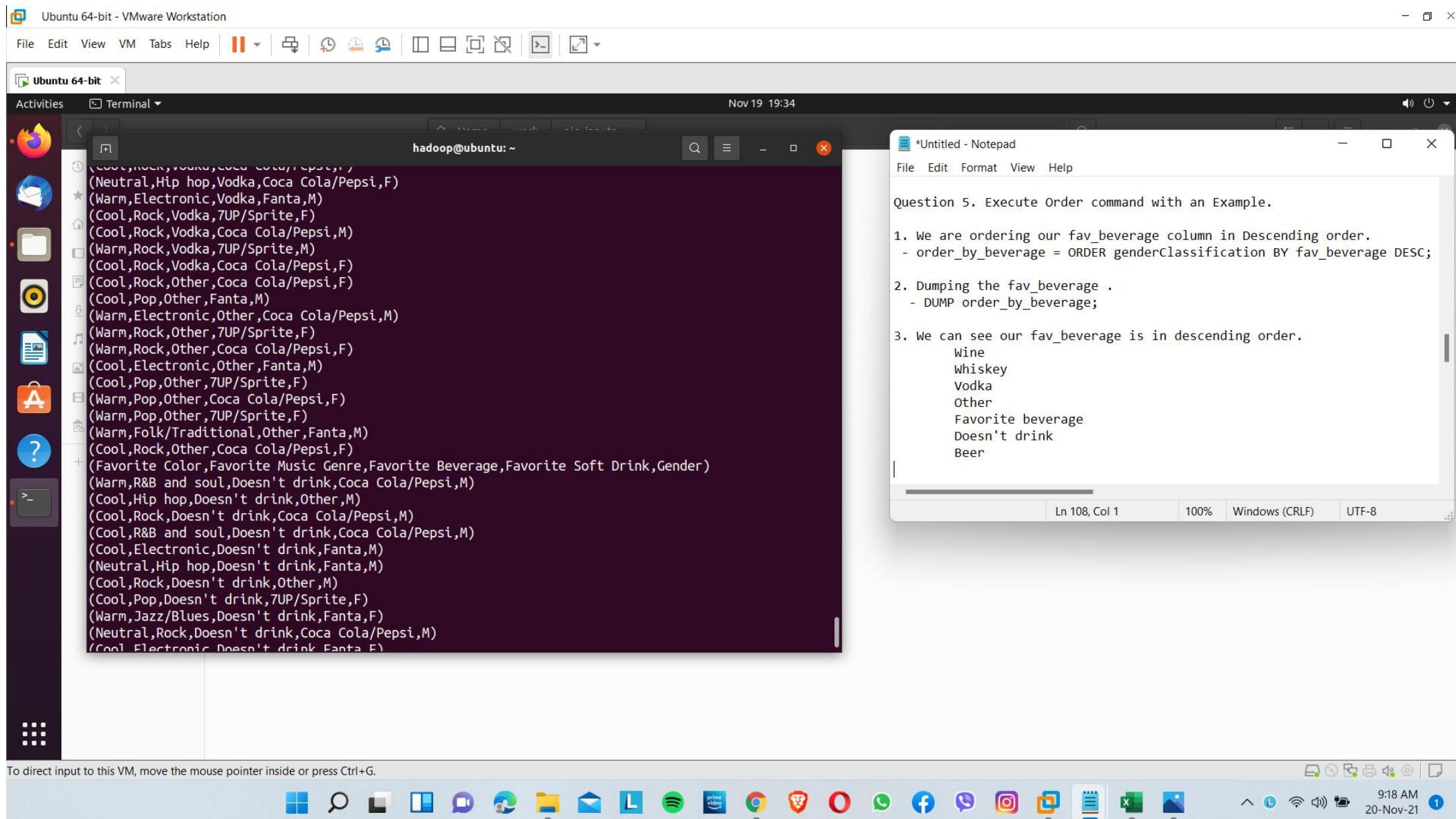


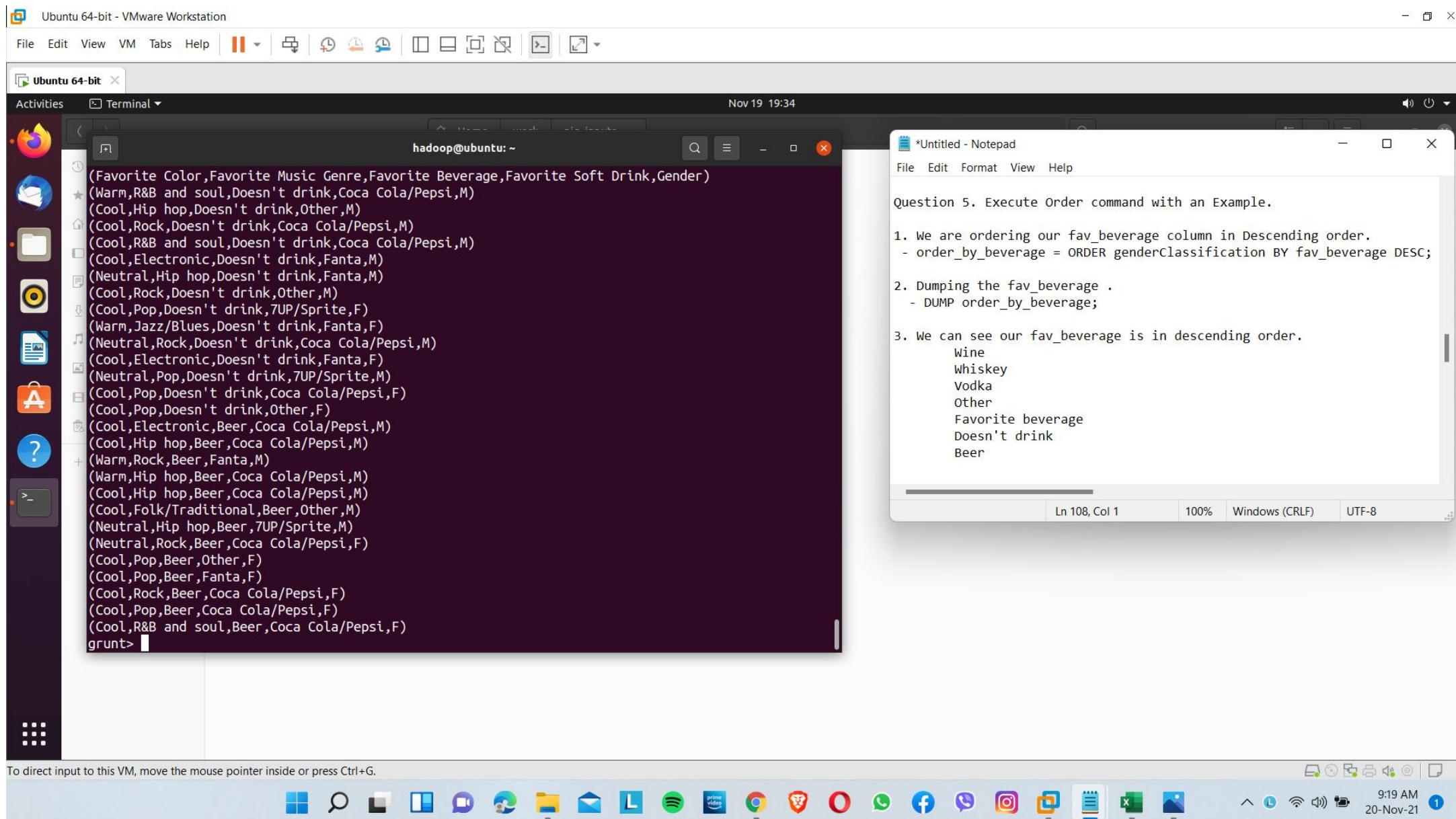












File Edit View VM Tabs Help | || | | | | | | | |

Ubuntu 64-bit

Activities Terminal

Nov 19 19:41

```
hadoop@ubuntu:~$ grunt> group_by_music_genre = GROUP genderClassification BY fav_music_genre;
grunt> DUMP group_by_music_genre;
2021-11-19 19:40:22,469 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used
in the script: GROUP_BY
2021-11-19 19:40:22,485 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default
.name is deprecated. Instead, use fs.defaultFS
2021-11-19 19:40:22,488 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple]
] was not set... will not generate code.
2021-11-19 19:40:22,515 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer
- {RULES_ENABLED=[AddForEach, ColumnMapKeyPrune, ConstantCalculator, GroupByConstParallelSetter,
LimitOptimizer, LoadTypeCastInserter, MergeFilter, MergeForEach, PartitionFilterOptimizer, Predica
tePushdownOptimizer, PushDownForEachFlatten, PushUpFilter, SplitFilter, StreamTypeCastInserter]}
2021-11-19 19:40:22,713 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.
MRCompiler - File concatenation threshold: 100 optimistic? false
2021-11-19 19:40:22,730 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.
MultiQueryOptimizer - MR plan size before optimization: 1
2021-11-19 19:40:22,730 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.
MultiQueryOptimizer - MR plan size after optimization: 1
2021-11-19 19:40:22,746 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default
.name is deprecated. Instead, use fs.defaultFS
2021-11-19 19:40:22,758 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - session.id
is deprecated. Instead, use dfs.metrics.session-id
2021-11-19 19:40:22,759 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Initializing JVM M
etrics with processName=JobTracker, sessionId=
2021-11-19 19:40:22,774 [main] INFO org.apache.pig.tools.pigstats.mapreduce.MRScriptState - Pig s
cript settings are added to the job
2021-11-19 19:40:22,778 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job
.reduce.markreset.buffer.percent is deprecated. Instead, use mapreduce.reduce.markreset.buffer.per
cent
```

*Untitled - Notepad

File Edit Format View Help

Question 6. Execute Group command with an Example.

1. We are grouping our data based on favorite music genre.
- group_by_music_genre = GROUP genderClassification BY fav_music_genre;
2. After that we are dumping our results.|

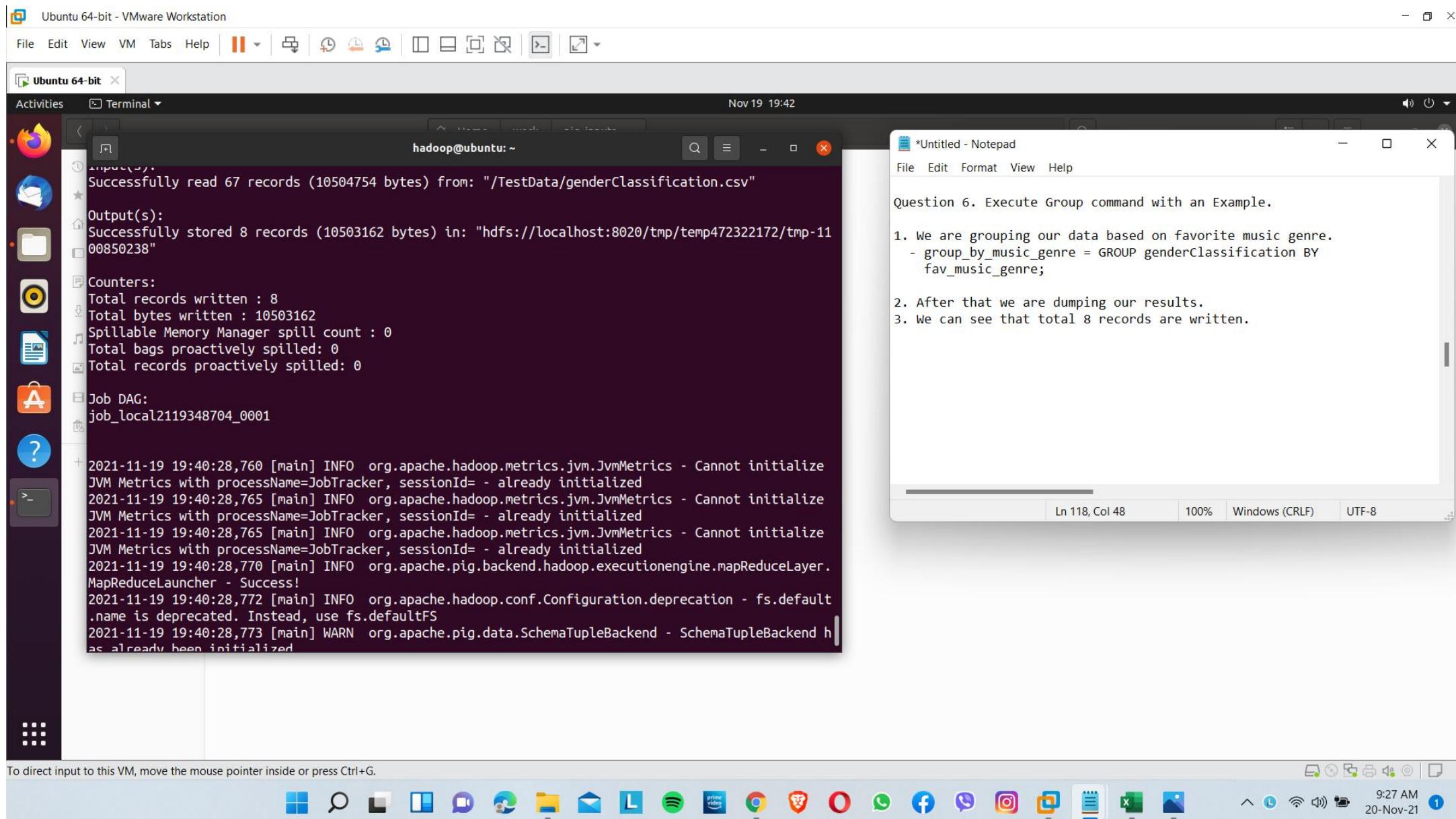
Ln 117, Col 42

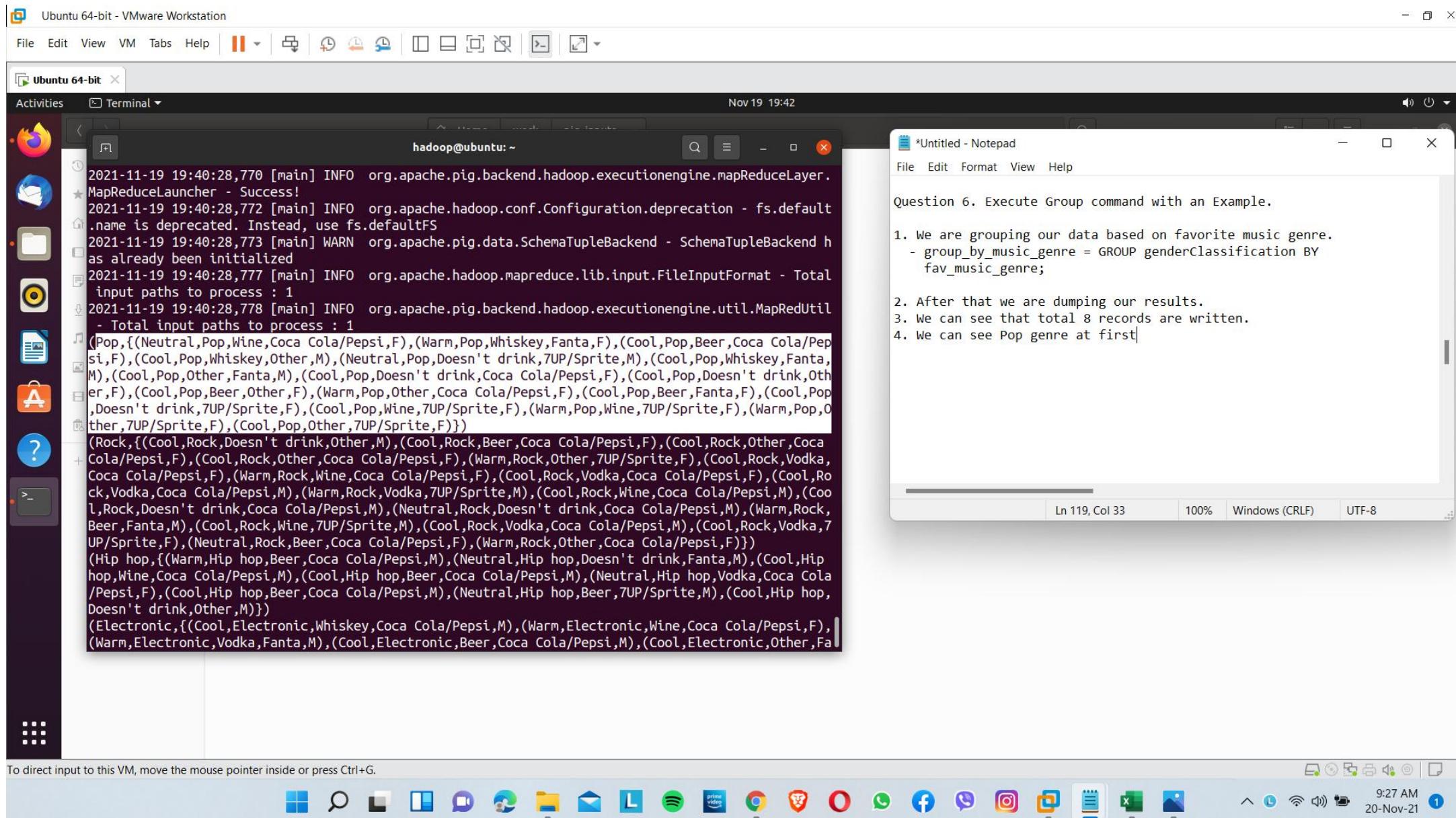
100% Windows (CRLF) UTF-8

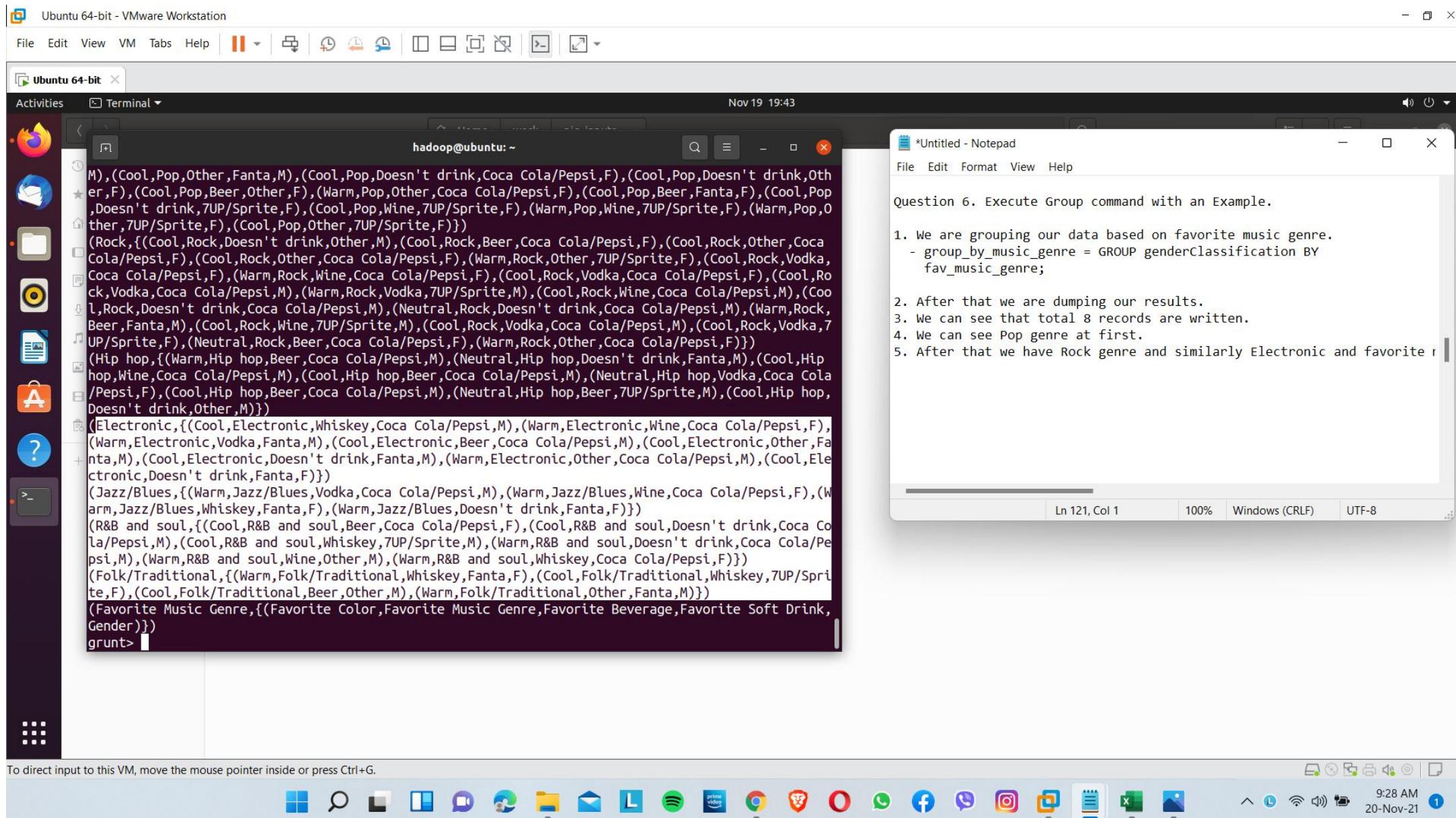
To direct input to this VM, move the mouse pointer inside or press Ctrl+G.

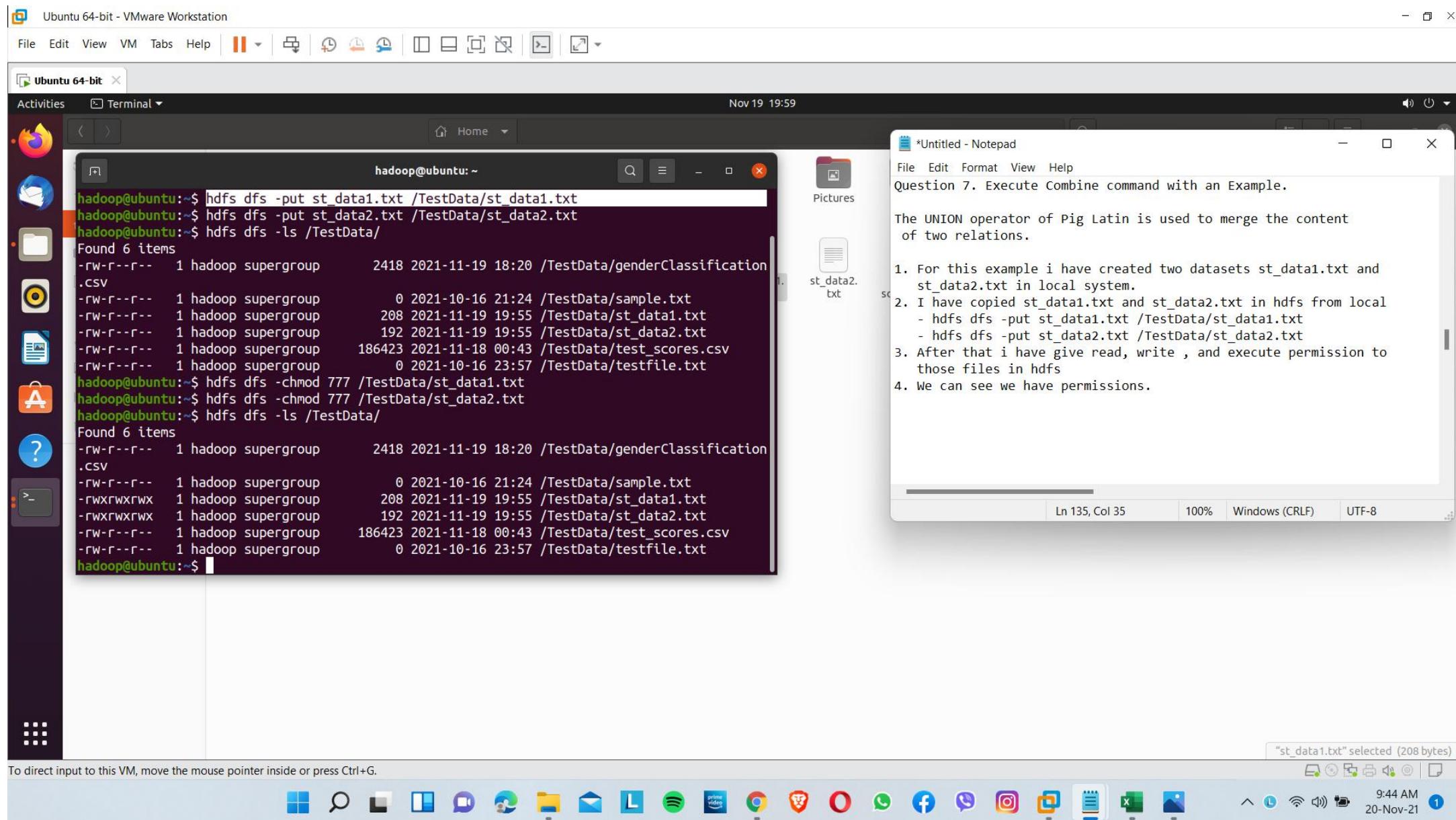


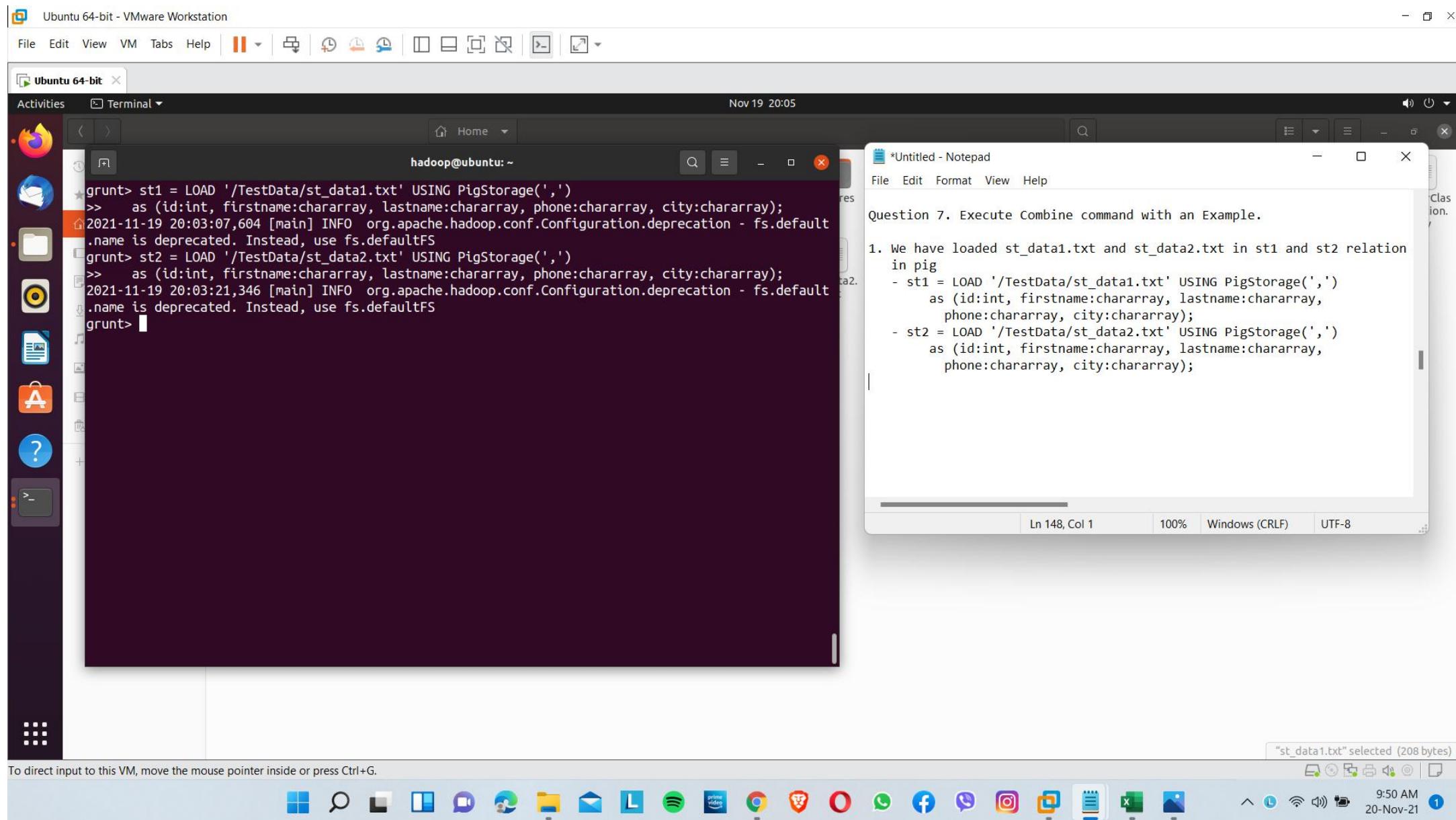
9:26 AM 20-Nov-21 1

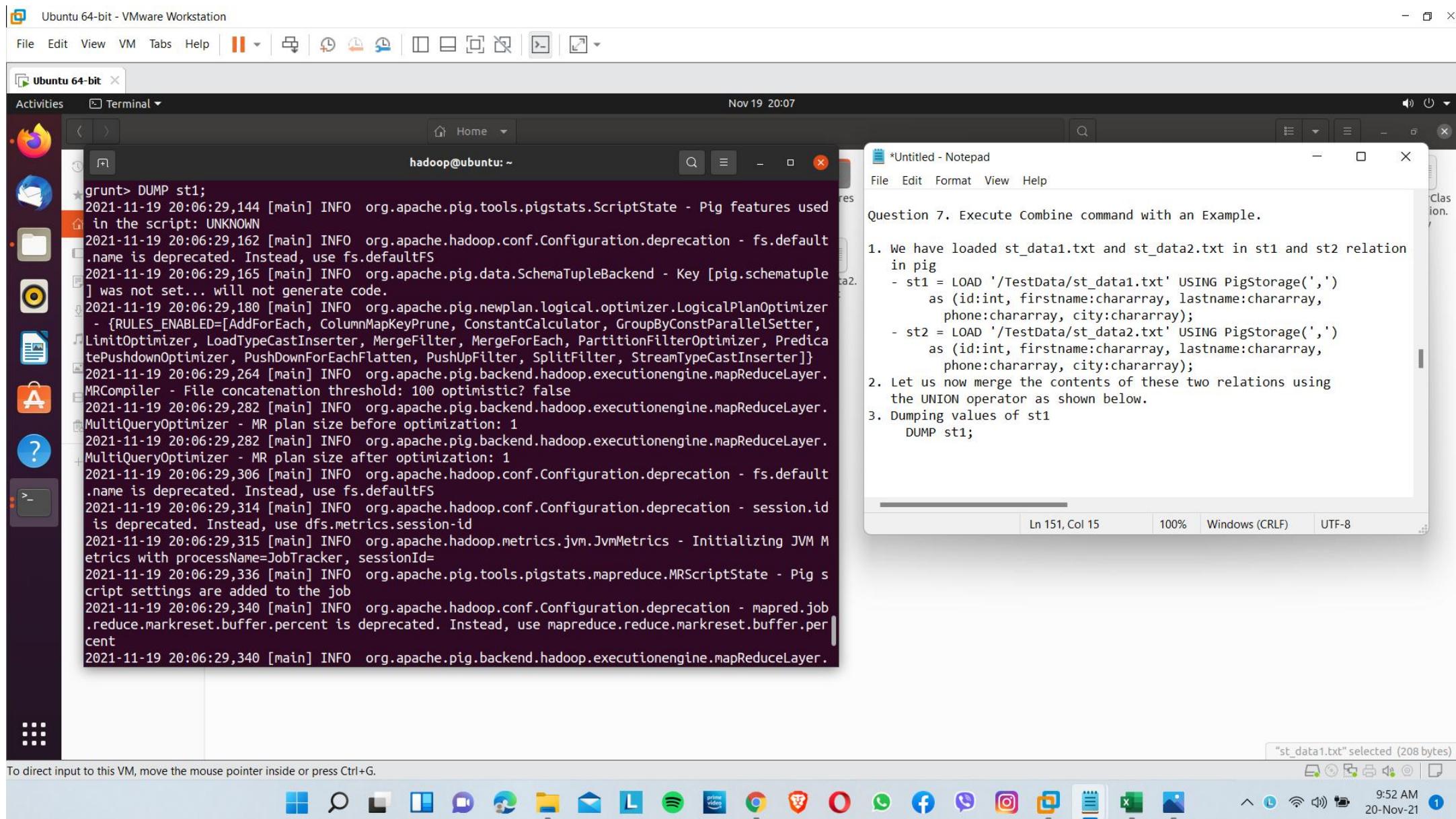












File Edit View VM Tabs Help | || | ⌂ ⌃ ⌄ ⌅ ⌆ ⌇ ⌈ ⌉ ⌊ ⌋

Ubuntu 64-bit

Activities Terminal

Nov 19 20:07

Clas
sion.

```
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_local1831513605_0001

2021-11-19 20:06:35,277 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize
JVM Metrics with processName=JobTracker, sessionId= - already initialized
2021-11-19 20:06:35,278 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize
JVM Metrics with processName=JobTracker, sessionId= - already initialized
2021-11-19 20:06:35,279 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize
JVM Metrics with processName=JobTracker, sessionId= - already initialized
2021-11-19 20:06:35,287 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.
MapReduceLauncher - Success!
2021-11-19 20:06:35,289 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default
.name is deprecated. Instead, use fs.defaultFS
2021-11-19 20:06:35,289 [main] WARN org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend h
as already been initialized
2021-11-19 20:06:35,295 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total
input paths to process : 1
2021-11-19 20:06:35,295 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil
- Total input paths to process : 1
(1,Mapreduce,Reddy,21,9848022337)
(2,Pig,Battacharya,22,9848022338)
(3,Hive,Khanna,22,9848022339)
(4,Flume,Agarwal,21,9848022330)
(5,Kafka,Mohanthy,23,9848022336)
grunt> |
```

*Untitled - Notepad

File Edit Format View Help

Question 7. Execute Combine command with an Example.

1. We have loaded st_data1.txt and st_data2.txt in st1 and st2 relation in pig
 - st1 = LOAD '/TestData/st_data1.txt' USING PigStorage(',')
as (id:int, firstname:chararray, lastname:chararray,
phone:chararray, city:chararray);
 - st2 = LOAD '/TestData/st_data2.txt' USING PigStorage(',')
as (id:int, firstname:chararray, lastname:chararray,
phone:chararray, city:chararray);
2. Let us now merge the contents of these two relations using the UNION operator as shown below.
3. Dumping values of st1
DUMP st1;

Ln 151, Col 15 100% Windows (CRLF) UTF-8

To direct input to this VM, move the mouse pointer inside or press Ctrl+G.

"st_data1.txt" selected (208 bytes)



9:52 AM
20-Nov-21 1

File Edit View VM Tabs Help | || | ⌂ ⌃ ⌄ ⌅ ⌆ ⌇ ⌈ ⌉ ⌊ ⌋

Ubuntu 64-bit

Activities Terminal

Nov 19 20:08

Clas
sion.

```
Total bags proactively spilled: 0
Total records proactively spilled: 0

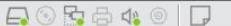
Job DAG:
job_local1433157370_0002

2021-11-19 20:07:50,416 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize
JVM Metrics with processName=JobTracker, sessionId= - already initialized
2021-11-19 20:07:50,416 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize
JVM Metrics with processName=JobTracker, sessionId= - already initialized
2021-11-19 20:07:50,417 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize
JVM Metrics with processName=JobTracker, sessionId= - already initialized
2021-11-19 20:07:50,418 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.
MapReduceLauncher - Success!
2021-11-19 20:07:50,419 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default
.name is deprecated. Instead, use fs.defaultFS
2021-11-19 20:07:50,419 [main] WARN org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend h
as already been initialized
2021-11-19 20:07:50,422 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total
input paths to process : 1
2021-11-19 20:07:50,423 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil
- Total input paths to process : 1
(6,Spark,Mishra,9848022335,Chennai)
(7,hadoop,Nayak,9848022334,trivendram)
(8,James,Nambiar,9848022333,Chennai)
(9,Python, Kumar, 9810937393, Bangalore)
(10, Java, Script, 009987283, Coding)
grunt> |
```

To direct input to this VM, move the mouse pointer inside or press Ctrl+G.



"st_data1.txt" selected (208 bytes)



9:53 AM
20-Nov-21 1

File Edit View VM Tabs Help | || | | | | | | |

Ubuntu 64-bit

Activities Terminal

Nov 19 20:10

```
hadoop@ubuntu: ~
grunt> st = UNION st1,st2 ;
grunt> DUMP st;
2021-11-19 20:08:47,474 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used
in the script: UNION
2021-11-19 20:08:47,483 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default
.name is deprecated. Instead, use fs.defaultFS
2021-11-19 20:08:47,484 [main] WARN org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend h
as already been initialized
2021-11-19 20:08:47,484 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer
- {RULES_ENABLED=[AddForEach, ColumnMapKeyPrune, ConstantCalculator, GroupByConstParallelSetter,
LimitOptimizer, LoadTypeCastInserter, MergeFilter, MergeForEach, PartitionFilterOptimizer, Predica
tePushdownOptimizer, PushDownForEachFlatten, PushUpFilter, SplitFilter, StreamTypeCastInserter]}
2021-11-19 20:08:47,486 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.
MRCompiler - File concatenation threshold: 100 optimistic? false
2021-11-19 20:08:47,487 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.
MultiQueryOptimizer - MR plan size before optimization: 1
2021-11-19 20:08:47,487 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.
MultiQueryOptimizer - MR plan size after optimization: 1
2021-11-19 20:08:47,492 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default
.name is deprecated. Instead, use fs.defaultFS
2021-11-19 20:08:47,493 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize
JVM Metrics with processName=JobTracker, sessionId= - already initialized
2021-11-19 20:08:47,494 [main] INFO org.apache.pig.tools.pigstats.mapreduce.MRScriptState - Pig s
cript settings are added to the job
2021-11-19 20:08:47,494 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.
JobControlCompiler - mapred.job.reduce.markreset.buffer.percent is not set, set to default 0.3
2021-11-19 20:08:47,494 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.
JobControlCompiler - This job cannot be converted run in-process
```

*Untitled - Notepad

File Edit Format View Help
Question 7. Execute Combine command with an Example.

1. We have loaded st_data1.txt and st_data2.txt in st1 and st2 relation in pig
 - st1 = LOAD '/TestData/st_data1.txt' USING PigStorage(',')
as (id:int, firstname:chararray, lastname:chararray,
phone:chararray, city:chararray);
 - st2 = LOAD '/TestData/st_data2.txt' USING PigStorage(',')
as (id:int, firstname:chararray, lastname:chararray,
phone:chararray, city:chararray);
2. Let us now merge the contents of these two relations using the UNION operator as shown below.
3. Dumping values of st1
 - DUMP st1;
4. Dumping st2 values
 - DUMP st2;
5. We are combining our but st1 and st2 file into single st file
 - UNION st1, st2;
 - DUMP st;

Ln 157, Col 13 100% Windows (CRLF) UTF-8

"st_data1.txt" selected (208 bytes)

To direct input to this VM, move the mouse pointer inside or press Ctrl+G.



9:55 AM 20-Nov-21 1

File Edit View VM Tabs Help | || | | | | | | |

Ubuntu 64-bit

Activities Terminal

Nov 19 20:10

Clas
sion.

```
hadoop@ubuntu: ~
2021-11-19 20:08:53,377 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize
JVM Metrics with processName=JobTracker, sessionId= - already initialized
2021-11-19 20:08:53,378 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize
JVM Metrics with processName=JobTracker, sessionId= - already initialized
2021-11-19 20:08:53,379 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize
JVM Metrics with processName=JobTracker, sessionId= - already initialized
2021-11-19 20:08:53,381 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.
MapReduceLauncher - Success!
2021-11-19 20:08:53,387 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default
.name is deprecated. Instead, use fs.defaultFS
2021-11-19 20:08:53,387 [main] WARN org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend h
as already been initialized
2021-11-19 20:08:53,396 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total
input paths to process : 2
2021-11-19 20:08:53,396 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil
- Total input paths to process : 2
(1,Mapreduce,Reddy,21,9848022337)
(2,Pig,Battacharya,22,9848022338)
(3,Hive,Khanna,22,9848022339)
(4,Flume,Agarwal,21,9848022330)
(5,Kafka,Mohanthy,23,9848022336)
(6,Spark,Mishra,9848022335,Chennai)
(7,hadoop,Nayak,9848022334,trivendram)
(8,James,Nambyayar,9848022333,Chennai)
(9,Python, Kumar, 9810937393, Bangalore)
(10, Java, Script, 009987283, Coding)
grunt> ■
```

To direct input to this VM, move the mouse pointer inside or press Ctrl+G.



"st_data1.txt" selected (208 bytes)



9:55 AM
20-Nov-21 1

*Untitled - Notepad

File Edit Format View Help
Question 7. Execute Combine command with an Example.

1. We have loaded st_data1.txt and st_data2.txt in st1 and st2 relation in pig
 - st1 = LOAD '/TestData/st_data1.txt' USING PigStorage(',')
as (id:int, firstname:chararray, lastname:chararray,
phone:chararray, city:chararray);
 - st2 = LOAD '/TestData/st_data2.txt' USING PigStorage(',')
as (id:int, firstname:chararray, lastname:chararray,
phone:chararray, city:chararray);
2. Let us now merge the contents of these two relations using the UNION operator as shown below.
3. Dumping values of st1
 - DUMP st1;
4. Dumping st2 values
 - DUMP st2;
5. We are combining our st1 and st2 file into single st file
 - UNION st1, st2;
 - DUMP st;
6. We can see all our data are combined together in a single file

File Edit View VM Tabs Help | || | | | | | | |

Ubuntu 64-bit

Activities Terminal

Nov 19 20:19

Speaker Power

```
hadoop@ubuntu:~$  
hadoop@ubuntu:~$ hdfs dfs -put customer.txt /TestData/customer.txt  
hadoop@ubuntu:~$ hdfs dfs -put order.txt /TestData/order.txt  
hadoop@ubuntu:~$ hdfs dfs -chmod 777 /TestData/customer.txt  
hadoop@ubuntu:~$ hdfs dfs -chmod 777 /TestData/order.txt  
hadoop@ubuntu:~$ hdfs dfs -ls /TestData/  
Found 8 items  
-rwxrwxrwx 1 hadoop supergroup 188 2021-11-19 20:15 /TestData/customer.txt  
-rw-r--r-- 1 hadoop supergroup 2418 2021-11-19 18:20 /TestData/genderClassification  
.csv  
-rwxrwxrwx 1 hadoop supergroup 125 2021-11-19 20:15 /TestData/order.txt  
-rw-r--r-- 1 hadoop supergroup 0 2021-10-16 21:24 /TestData/sample.txt  
-rwxrwxrwx 1 hadoop supergroup 208 2021-11-19 19:55 /TestData/st_data1.txt  
-rwxrwxrwx 1 hadoop supergroup 192 2021-11-19 19:55 /TestData/st_data2.txt  
not-rw-r--r-- 1 hadoop supergroup 186423 2021-11-18 00:43 /TestData/test_scores.csv  
-rw-r--r-- 1 hadoop supergroup 0 2021-10-16 23:57 /TestData/testfile.txt  
hadoop@ubuntu:~$
```

*Untitled - Notepad

File Edit Format View Help

9. Execute Inner Join command with an Example.

1. for this example , i have created two text file customer.txt and order.txt file which both has id column in common.
2. I have copied customer.txt and order.txt from local to hdfs
 - hdfs dfs -put customer.txt /TestData/customer.txt
 - hdfs dfs -put order.txt /TestData/order.txt
3. After that i gave all permission to both files in hdfs
 - dfs dfs -chmod 777 /TestData/customer.txt
 - dfs dfs -chmod 777 /TestData/order.txt

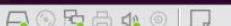
Ln 169, Col 39

100%

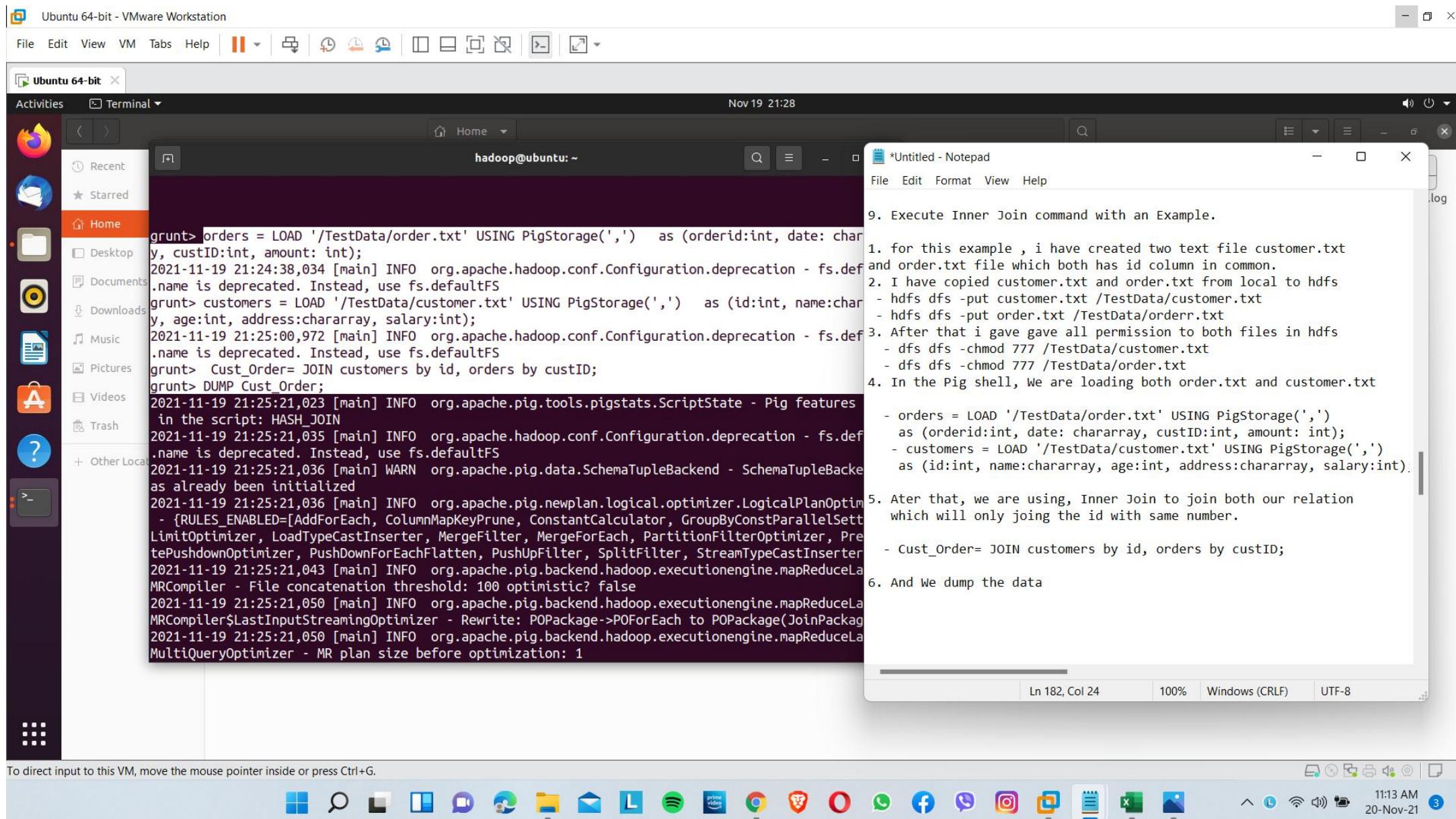
Windows (CRLF)

UTF-8

To direct input to this VM, move the mouse pointer inside or press Ctrl+G.



10:04 AM 20-Nov-21 1



File Edit View VM Help | | | | |

Ubuntu 64-bit

Activities Terminal Nov 19 21:30

```
hadoop@ubuntu:~$ Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_local1404424815_0003

2021-11-19 21:25:27,873 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize JVM Metrics with processName=JobTracker, sessionId= - already initialized
2021-11-19 21:25:27,873 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize JVM Metrics with processName=JobTracker, sessionId= - already initialized
2021-11-19 21:25:27,874 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize JVM Metrics with processName=JobTracker, sessionId= - already initialized
2021-11-19 21:25:27,877 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - MapReduceLauncher - Success!
2021-11-19 21:25:27,877 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.defaultFS .name is deprecated. Instead, use fs.defaultFS
2021-11-19 21:25:27,878 [main] WARN org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend as already been initialized
2021-11-19 21:25:27,880 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-11-19 21:25:27,880 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRed - Total input paths to process : 1
(2,Khilan,25,Delhi,1500,101,2009-11-20 00:00:00,2,1560)
(3,kaushik,23,Kota,2000,100,2009-10-08 00:00:00,3,1500)
(3,kaushik,23,Kota,2000,102,2009-10-08 00:00:00,3,3000)
(4,Chaitali,25,Mumbai,6500,103,2008-05-20 00:00:00,4,2060)
grunt> 
```

*Untitled - Notepad

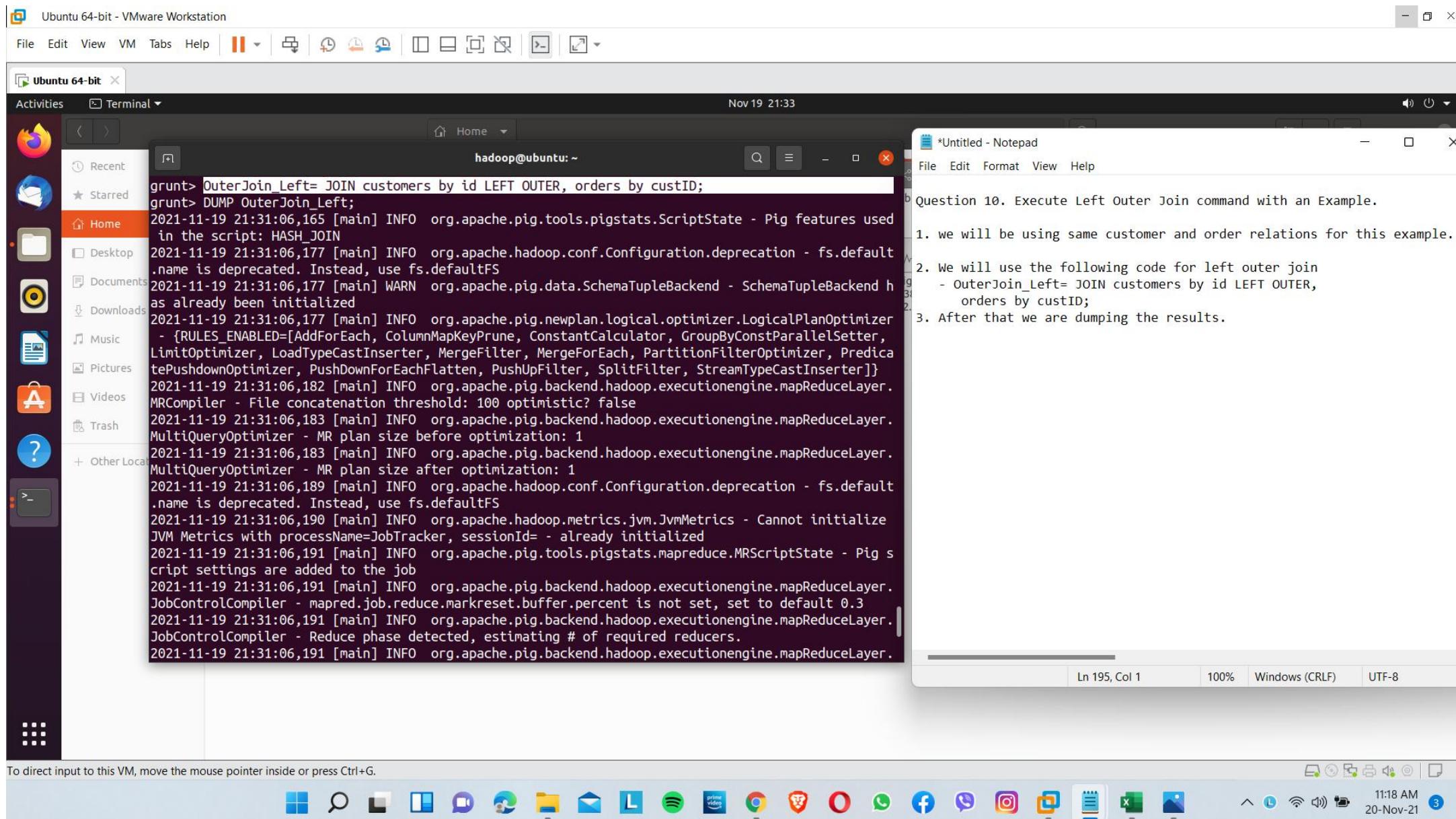
File Edit Format View Help

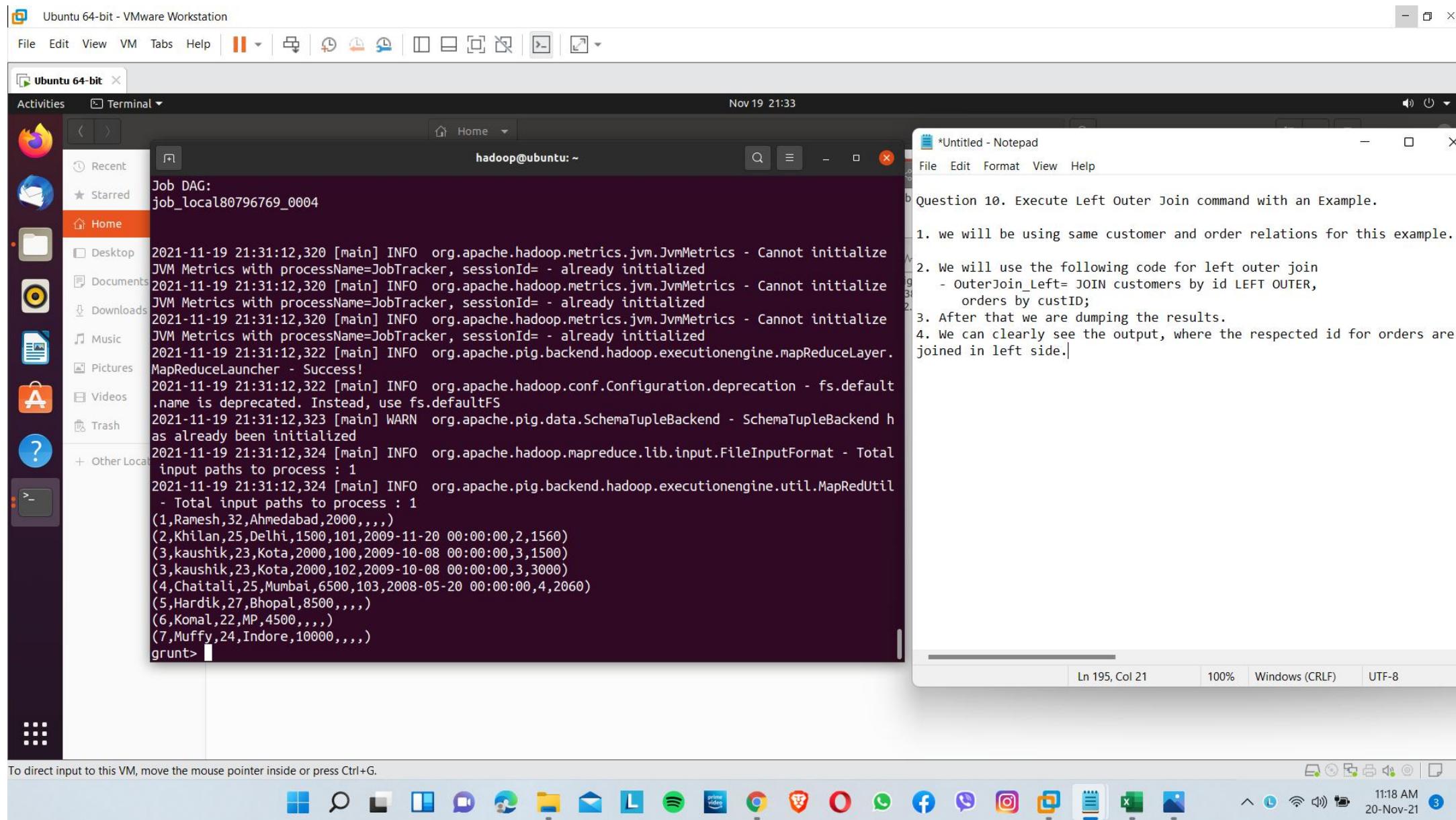
9. Execute Inner Join command with an Example.

- for this example , i have created two text file customer.txt and order.txt file which both has id column in common.
- I have copied customer.txt and order.txt from local to hdfs
 - hdfs dfs -put customer.txt /TestData/customer.txt
 - hdfs dfs -put order.txt /TestData/order.txt
- After that i gave all permission to both files in hdfs
 - dfs dfs -chmod 777 /TestData/customer.txt
 - dfs dfs -chmod 777 /TestData/order.txt
- In the Pig shell, We are loading both order.txt and customer.txt
 - orders = LOAD '/TestData/order.txt' USING PigStorage(',') as (orderid:int, date: chararray, custID:int, amount: int);
 - customers = LOAD '/TestData/customer.txt' USING PigStorage(',') as (id:int, name:chararray, age:int, address:chararray, salary:int);
- Ater that, we are using, Inner Join to join both our relation which will only joing the id with same number.
 - Cust_Order= JOIN customers by id, orders by custID;
- And We dump the data
- We can see clearly that only same id from both relations are printed.

Ln 184, Col 1 100% Windows (CRLF) UTF-8

To direct input to this VM, move the mouse pointer inside or press Ctrl+G.





File Edit View VM Help | II | | | | | | |

Ubuntu 64-bit X

Activities Terminal ▾

Nov 19 21:36

Speaker Power



Recent



Starred

Home



Desktop



Documents



Downloads



Music



Pictures



Videos



Trash



Other Local



```
(7,Muffy,24,Indore,10000,,,)  
grunt>  
grunt> OuterJoin_Right= JOIN customers by id RIGHT OUTER, orders by custID;  
grunt> DUMP OuterJoin_Right;  
2021-11-19 21:34:54,921 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used  
in the script: HASH_JOIN  
2021-11-19 21:34:54,931 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default  
.name is deprecated. Instead, use fs.defaultFS  
2021-11-19 21:34:54,932 [main] WARN org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend h  
as already been initialized  
2021-11-19 21:34:54,932 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer  
- {RULES_ENABLED=[AddForEach, ColumnMapKeyPrune, ConstantCalculator, GroupByConstParallelSetter,  
LimitOptimizer, LoadTypeCastInserter, MergeFilter, MergeForEach, PartitionFilterOptimizer, Predica  
tePushdownOptimizer, PushDownForEachFlatten, PushUpFilter, SplitFilter, StreamTypeCastInserter]}  
2021-11-19 21:34:54,934 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.  
MRCompiler - File concatenation threshold: 100 optimistic? false  
2021-11-19 21:34:54,935 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.  
MultiQueryOptimizer - MR plan size before optimization: 1  
2021-11-19 21:34:54,935 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.  
MultiQueryOptimizer - MR plan size after optimization: 1  
2021-11-19 21:34:54,941 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default  
.name is deprecated. Instead, use fs.defaultFS  
2021-11-19 21:34:54,942 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize  
JVM Metrics with processName=JobTracker, sessionId= - already initialized  
2021-11-19 21:34:54,942 [main] INFO org.apache.pig.tools.pigstats.mapreduce.MRScriptState - Pig s  
cript settings are added to the job  
2021-11-19 21:34:54,943 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.  
JobControlCompiler - mapred.job.reduce.markreset.buffer.percent is not set, set to default 0.3  
2021-11-19 21:34:54,943 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.
```

*Untitled - Notepad

File Edit Format View Help

Question 11. Execute Right outer join command with an Example.

1. We will use the following code to join Right Outer Join
 - OuterJoin_Right= JOIN customers by id RIGHT OUTER, orders by custID;
2. After that we are dumping the results

Ln 204, Col 41

100%

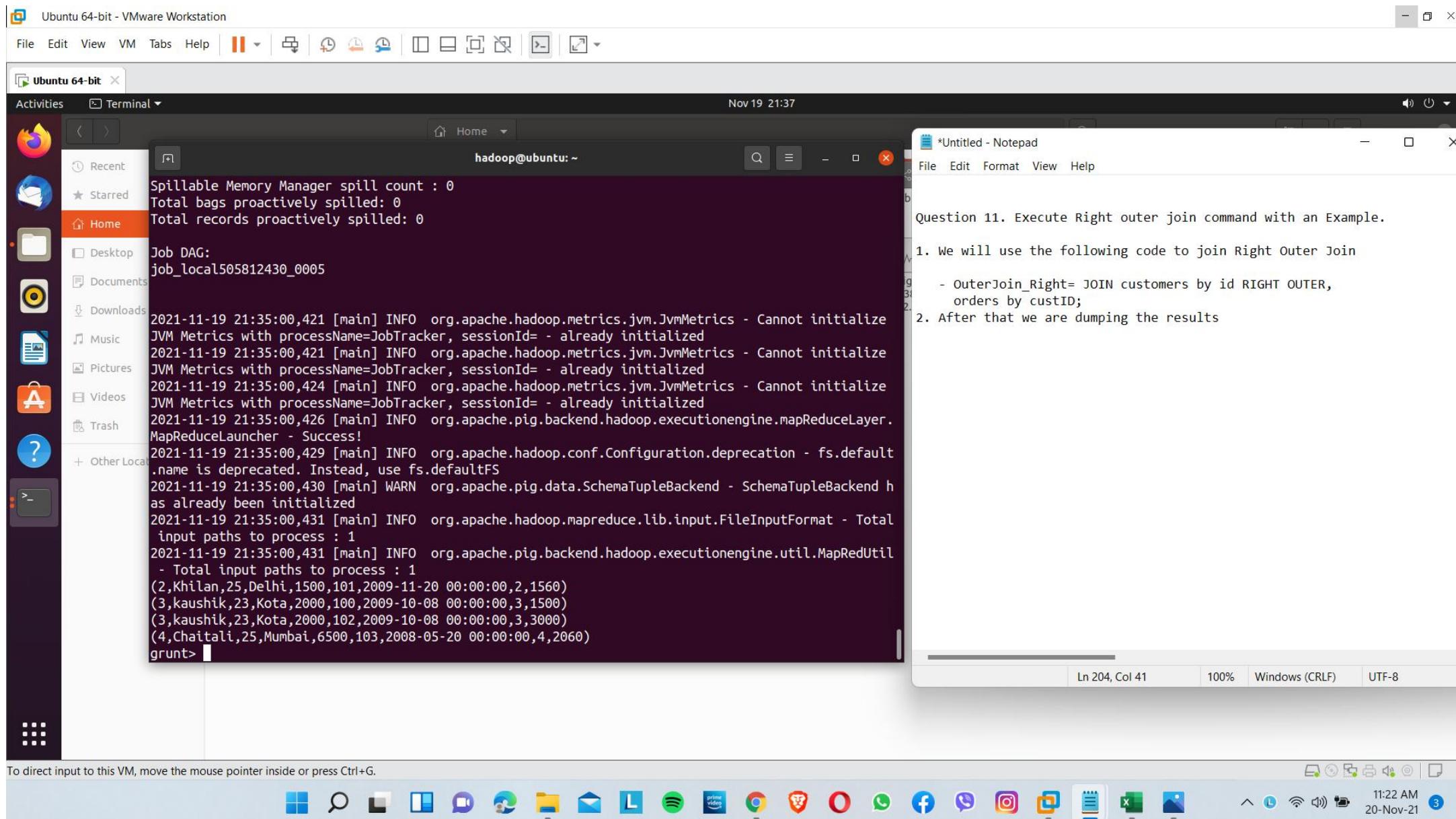
Windows (CRLF)

UTF-8

To direct input to this VM, move the mouse pointer inside or press Ctrl+G.



11:21 AM
20-Nov-21 3



File Edit View VM Help | || | | | | | | |

Ubuntu 64-bit

Activities Terminal

Nov 19 21:42

Speaker Power



Recent
Starred

Home

Desktop

Documents

Downloads

Music

Pictures

Videos

Trash

Other Local

```
grunt> FullOuterJoin= JOIN customers by id FULL OUTER, orders by custID;
grunt> DUMP FullOuterJoin;
2021-11-19 21:41:07,319 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used
in the script: HASH_JOIN
2021-11-19 21:41:07,331 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default
.name is deprecated. Instead, use fs.defaultFS
2021-11-19 21:41:07,331 [main] WARN org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend h
as already been initialized
2021-11-19 21:41:07,331 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer
- {RULES_ENABLED=[AddForEach, ColumnMapKeyPrune, ConstantCalculator, GroupByConstParallelSetter,
LimitOptimizer, LoadTypeCastInserter, MergeFilter, MergeForEach, PartitionFilterOptimizer, Predica
tePushdownOptimizer, PushDownForEachFlatten, PushUpFilter, SplitFilter, StreamTypeCastInserter]}
2021-11-19 21:41:07,333 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.
MRCompiler - File concatenation threshold: 100 optimistic? false
2021-11-19 21:41:07,334 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.
MultiQueryOptimizer - MR plan size before optimization: 1
2021-11-19 21:41:07,334 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.
MultiQueryOptimizer - MR plan size after optimization: 1
2021-11-19 21:41:07,339 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default
.name is deprecated. Instead, use fs.defaultFS
2021-11-19 21:41:07,341 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize
JVM Metrics with processName=JobTracker, sessionId= - already initialized
2021-11-19 21:41:07,341 [main] INFO org.apache.pig.tools.pigstats.mapreduce.MRScriptState - Pig s
cript settings are added to the job
2021-11-19 21:41:07,341 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.
JobControlCompiler - mapred.job.reduce.markreset.buffer.percent is not set, set to default 0.3
2021-11-19 21:41:07,341 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.
JobControlCompiler - Reduce phase detected, estimating # of required reducers.
2021-11-19 21:41:07,342 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.
```

*Untitled - Notepad

File Edit Format View Help

Question 12. Execute Full Join command with an Example.

1. We will use the following code to join Right Outer Join

- FullOuterJoin= JOIN customers by id FULL OUTER,
orders by custID;

2. After that we are dumping the results
-DUMP FullOuterJoin

Ln 216, Col 25

100%

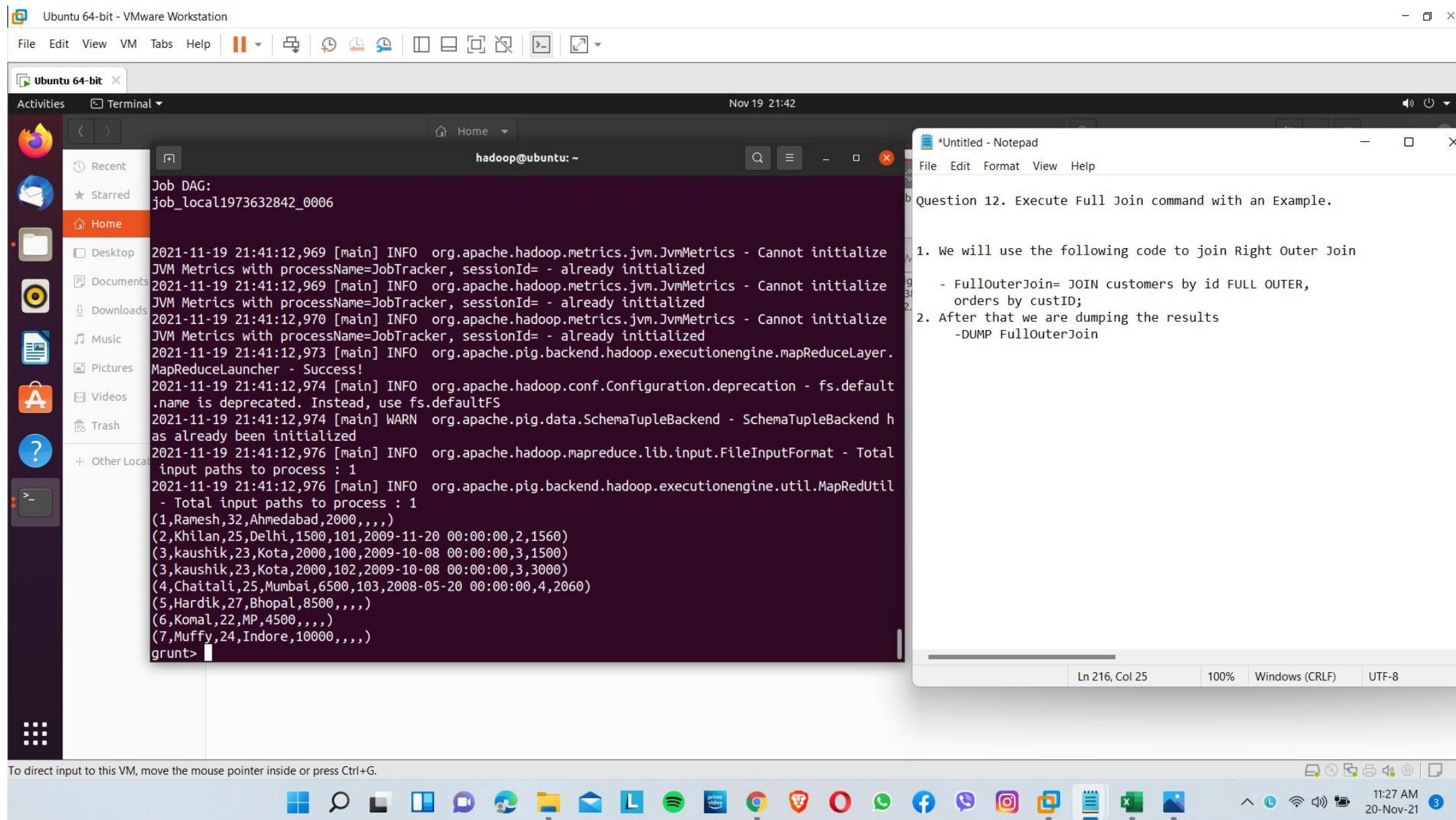
Windows (CRLF)

UTF-8

To direct input to this VM, move the mouse pointer inside or press Ctrl+G.

File Explorer Search Task View Taskbar

11:27 AM 20-Nov-21 3



File Edit View VM Help | || | | | | | | |

Ubuntu 64-bit

Activities Terminal

Nov 19 22:17

Speaker Power

```
(6,Komal,22,MP,4500,,,)  
(7,Muffy,24,Indore,10000,,,)  
=====grunt> cogroup_data = COGROUP customers by id, orders by custID;  
hadoop@grunt> DUMP cogroup_data;  
hadoop@2021-11-19 22:15:05,636 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used  
ls: `/Te in the script: COGROUP  
hadoop@2021-11-19 22:15:05,648 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default  
ls: `/Te.name is deprecated. Instead, use fs.defaultFS  
hadoop@2021-11-19 22:15:05,649 [main] WARN org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend h  
Found 8 as already been initialized  
-rwxrwxr 2021-11-19 22:15:05,649 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer  
-rw-r--r - {RULES_ENABLED=[AddForEach, ColumnMapKeyPrune, ConstantCalculator, GroupByConstParallelSetter,  
.csv LimitOptimizer, LoadTypeCastInserter, MergeFilter, MergeForEach, PartitionFilterOptimizer, Predica  
-rwxrwxr tePushdownOptimizer, PushDownForEachFlatten, PushUpFilter, SplitFilter, StreamTypeCastInserter]}  
-rw-r--r 2021-11-19 22:15:05,658 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.  
-rwxrwxr MRCompiler - File concatenation threshold: 100 optimistic? false  
-rwxrwxr 2021-11-19 22:15:05,659 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.  
-rw-r--r MultiQueryOptimizer - MR plan size before optimization: 1  
-rw-r--r 2021-11-19 22:15:05,659 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.  
hadoop@MultiQueryOptimizer - MR plan size after optimization: 1  
hadoop@2021-11-19 22:15:05,664 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default  
hadoop@.name is deprecated. Instead, use fs.defaultFS  
2021-11-19 22:15:05,664 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize  
JVM Metrics with processName=JobTracker, sessionId= - already initialized  
2021-11-19 22:15:05,665 [main] INFO org.apache.pig.tools.pigstats.mapreduce.MRScriptState - Pig s  
cript settings are added to the job  
2021-11-19 22:15:05,666 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.  
JobControlCompiler - mapred.job.reduce.markreset.buffer.percent is not set, set to default 0.3  
2021-11-19 22:15:05,666 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.
```

*Untitled - Notepad

File Edit Format View Help

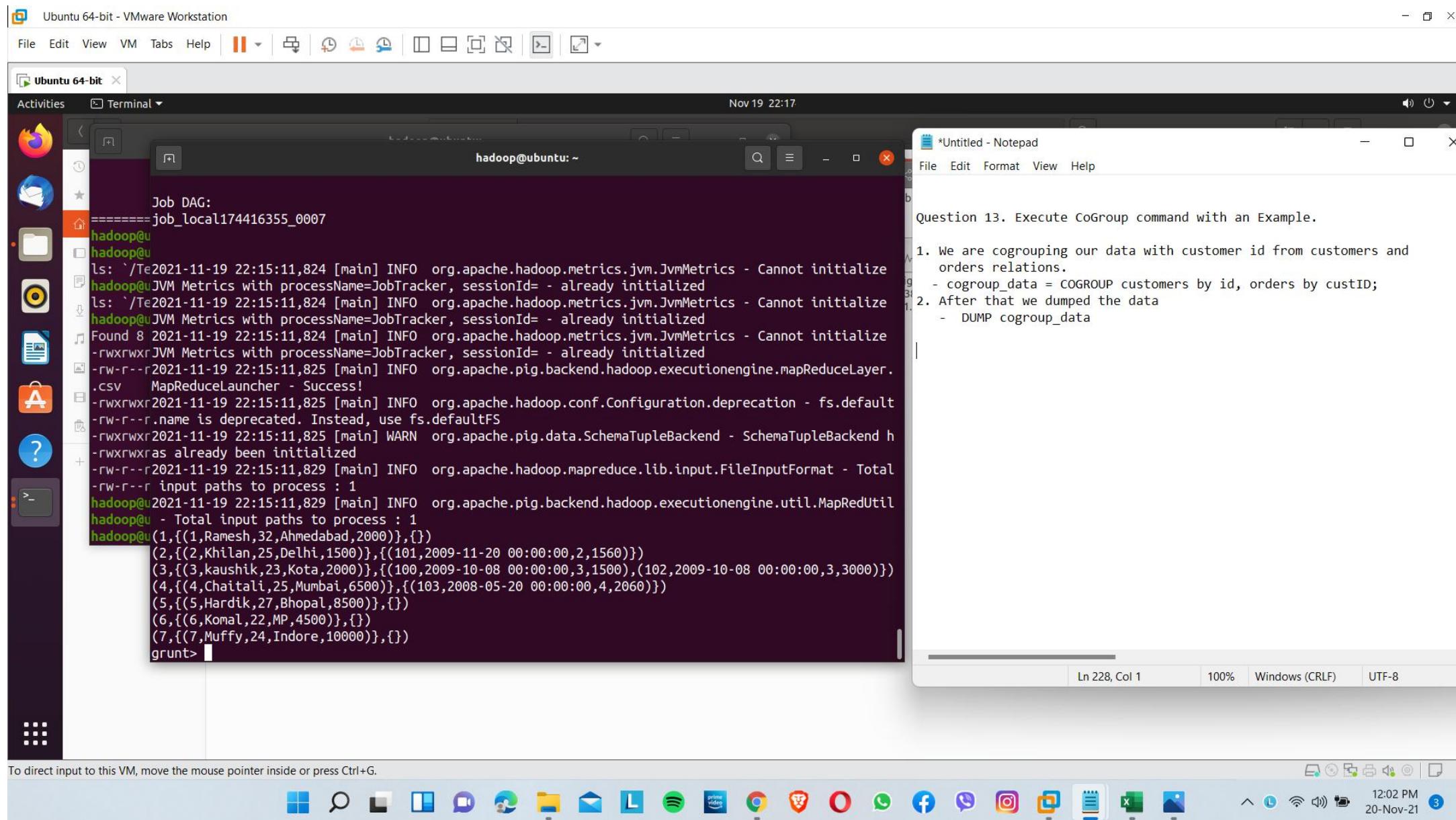
Question 13. Execute CoGroup command with an Example.

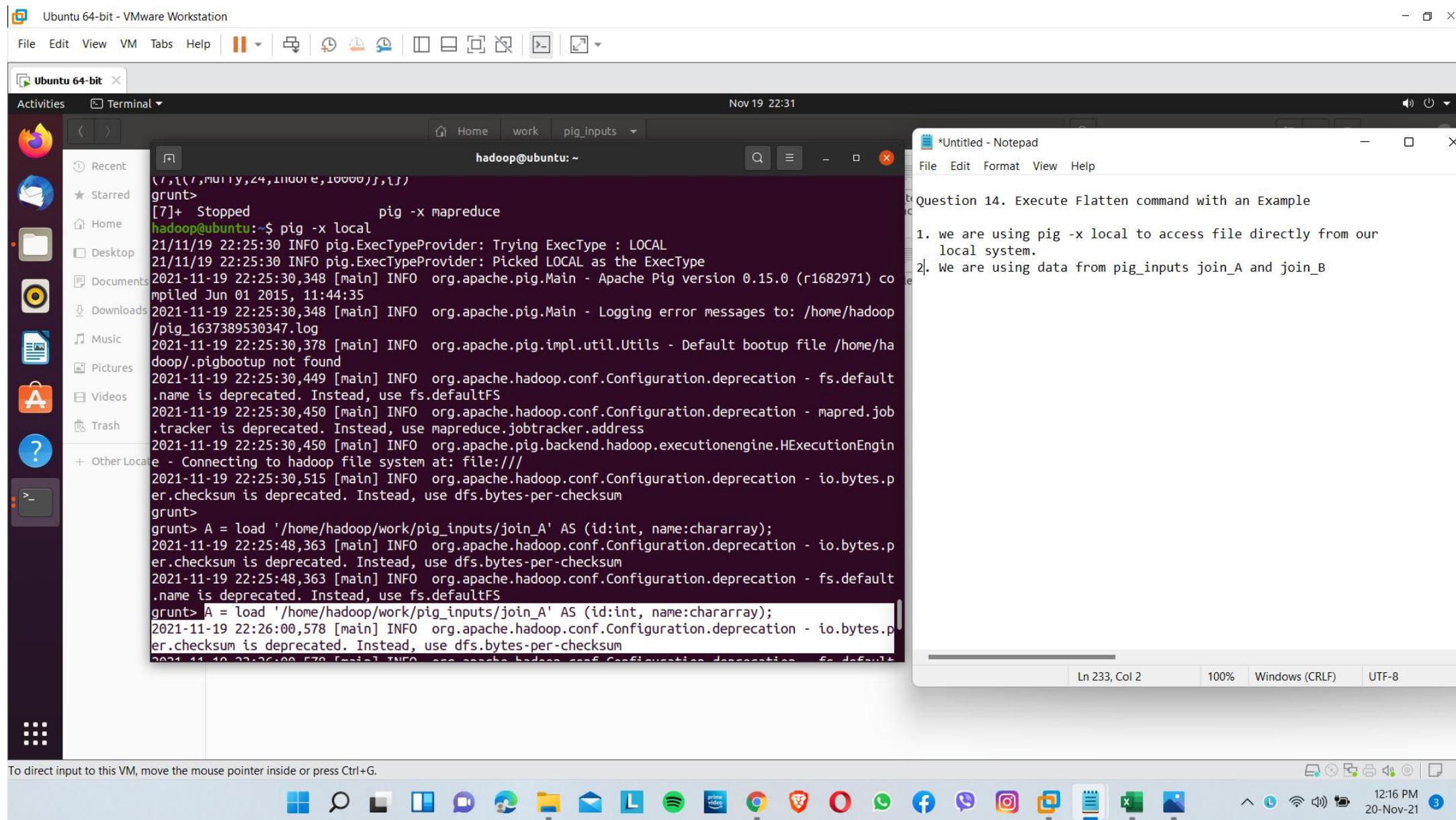
1. We are cogrouping our data with customer id from customers and orders relations.
 - cogroup_data = COGROUP customers by id, orders by custID;
2. After that we dumped the data
 - DUMP cogroup_data

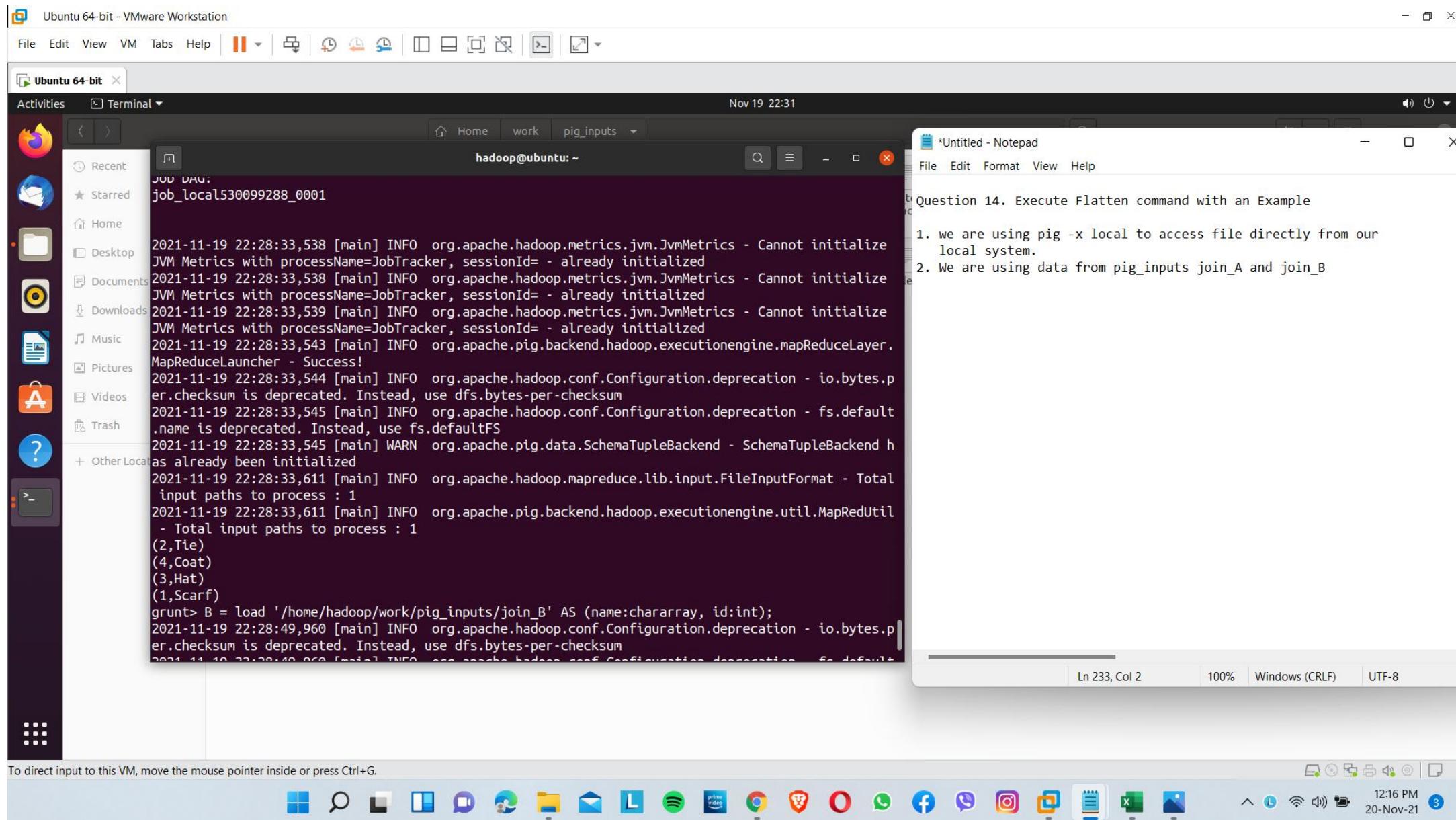
To direct input to this VM, move the mouse pointer inside or press Ctrl+G.

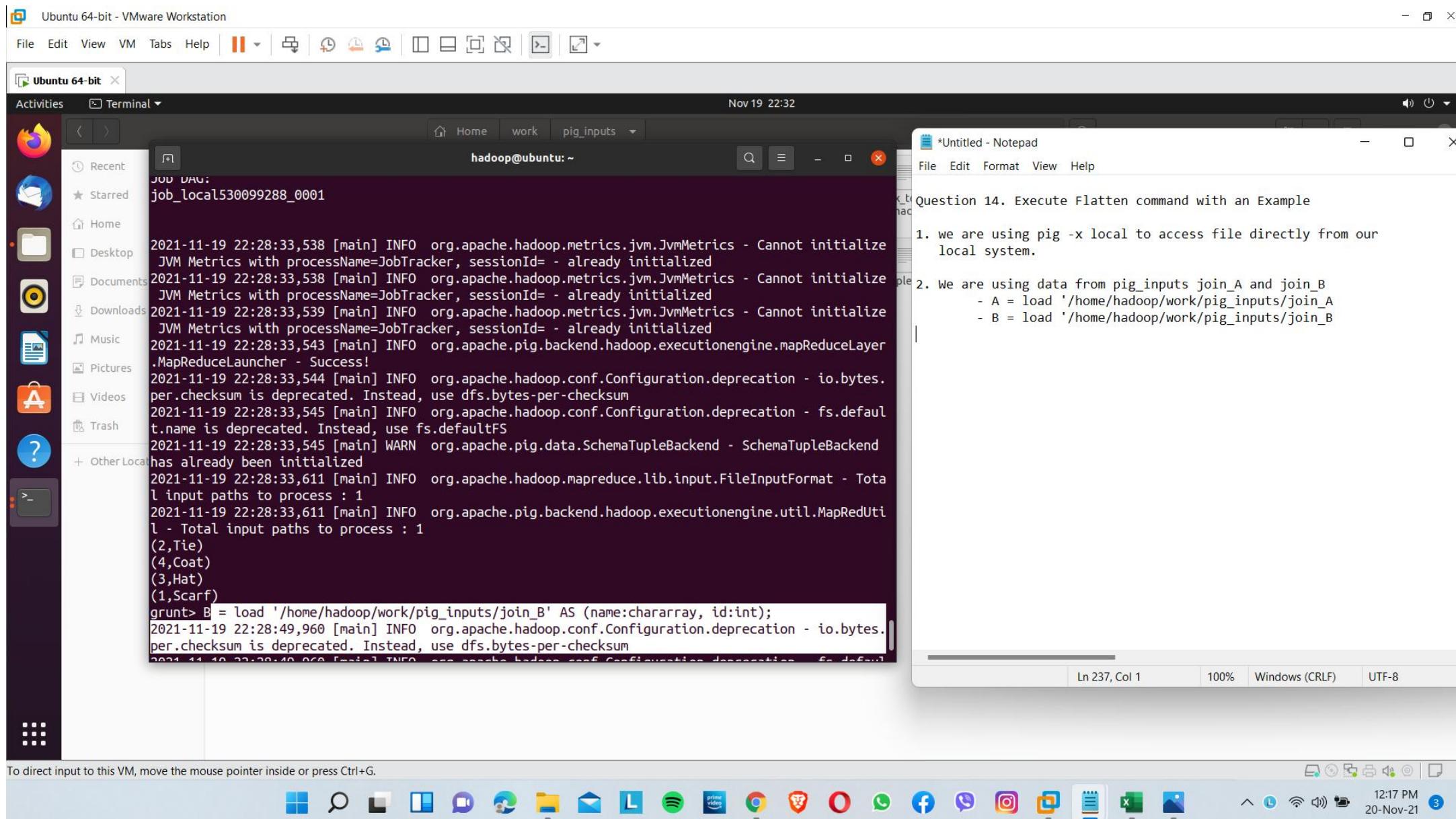


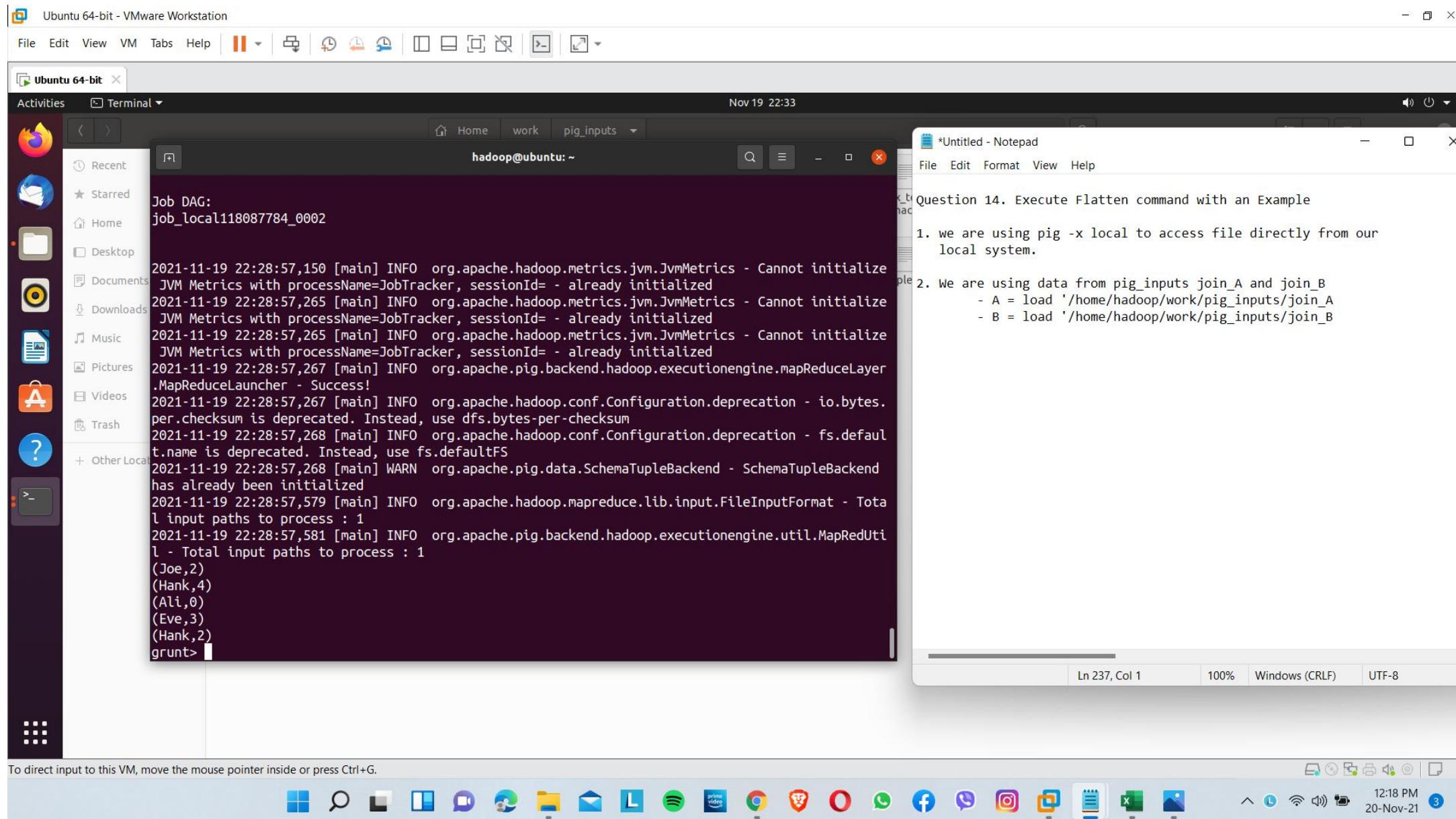
12:02 PM 20-Nov-21 3

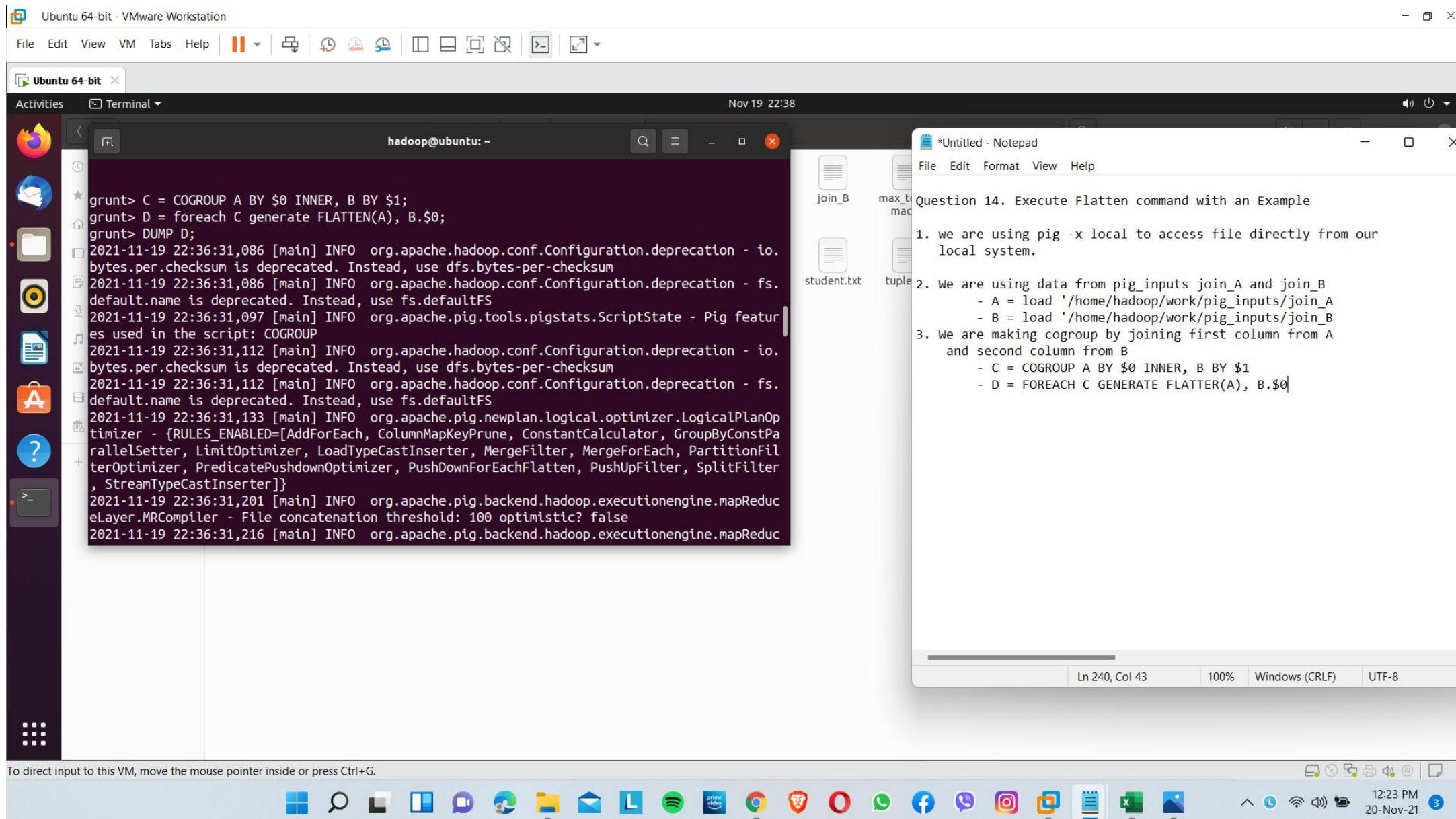


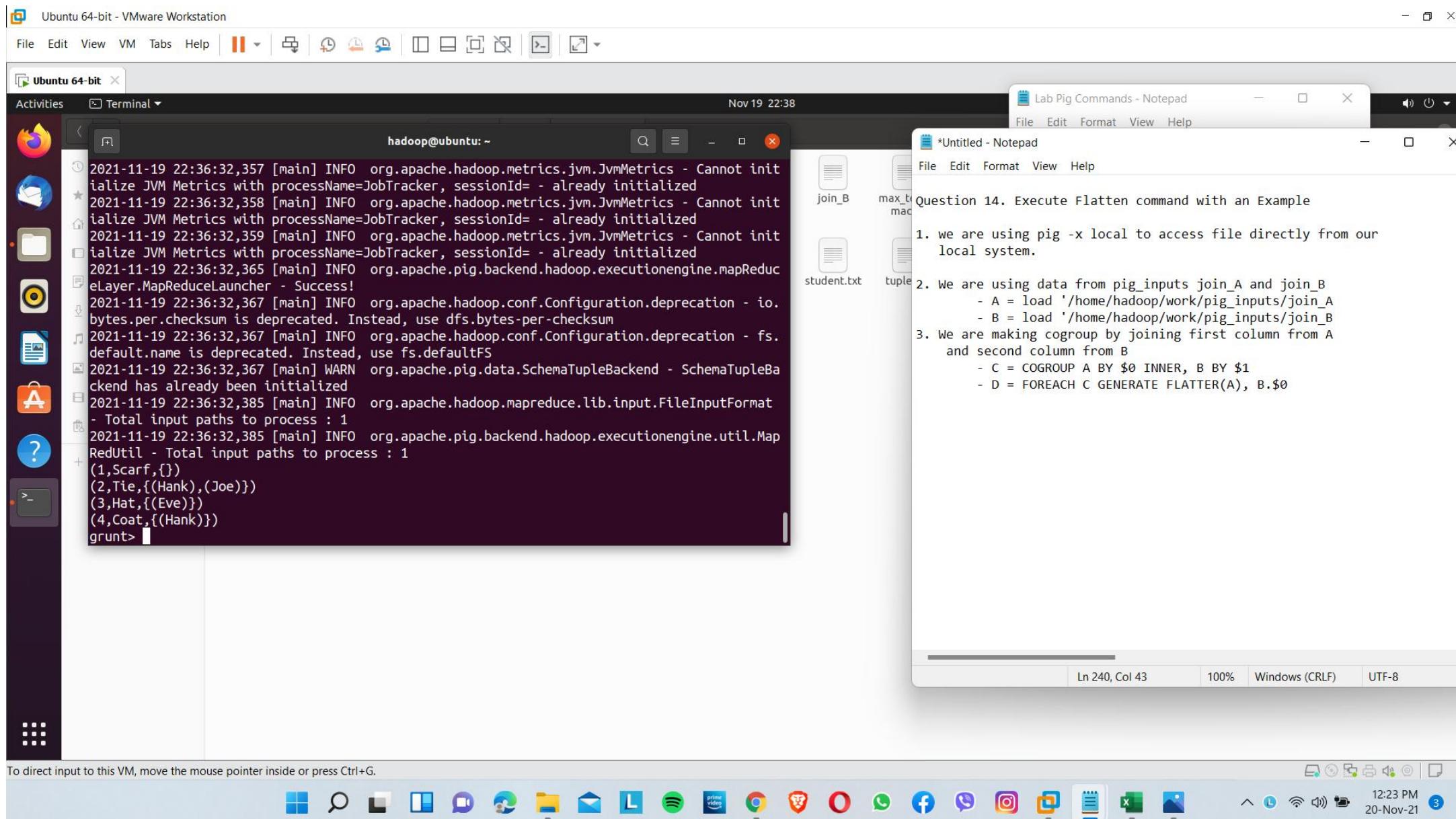












File Edit View VM Help | || | | | | | | | |

Ubuntu 64-bit

Activities Terminal

Nov 19 22:43

Speaker Power



hadoop@ubuntu: ~

Q E

```
hadoop@ubuntu:~$ pig
21/11/19 22:40:58 INFO pig.ExecTypeProvider: Trying ExecType : LOCAL
21/11/19 22:40:58 INFO pig.ExecTypeProvider: Trying ExecType : MAPREDUCE
21/11/19 22:40:58 INFO pig.ExecTypeProvider: Picked MAPREDUCE as the ExecType
2021-11-19 22:40:58,978 [main] INFO org.apache.pig.Main - Apache Pig version 0.15.0 (r1682971) compiled
01 2015, 11:44:35
2021-11-19 22:40:58,978 [main] INFO org.apache.pig.Main - Logging error messages to: /home/hadoop/pig_
0458977.log
2021-11-19 22:40:58,988 [main] INFO org.apache.pig.impl.util.Utils - Default bootup file /home/hadoop/
bootup not found
2021-11-19 22:40:59,181 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.trac
deprecated. Instead, use mapreduce.jobtracker.address
2021-11-19 22:40:59,181 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name
prected. Instead, use fs.defaultFS
2021-11-19 22:40:59,181 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - C
onnecting to hadoop file system at: hdfs://localhost:8020
2021-11-19 22:41:00,033 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name
prected. Instead, use fs.defaultFS
grunt>
```

```
grunt> customers = LOAD '/TestData/customer.txt' USING PigStorage(',');
>>   as (id:int, name:chararray, age:int, address:chararray, salary:int);
2021-11-19 22:41:39,755 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name
prected. Instead, use fs.defaultFS
grunt> tobag = FOREACH customers GENERATE TOBAG (id,name,age,address,salary);
grunt> DUMP tobag;
2021-11-19 22:42:32,751 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in t
his script: UNKNOWN
2021-11-19 22:42:32,773 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is de
```

*Untitled - Notepad

File Edit Format View Help

Question 15. Execute ToBag command with an Example.

For this example, we are using pig default mode

Ln 245, Col 1

100%

Windows (CRLF)

UTF-8

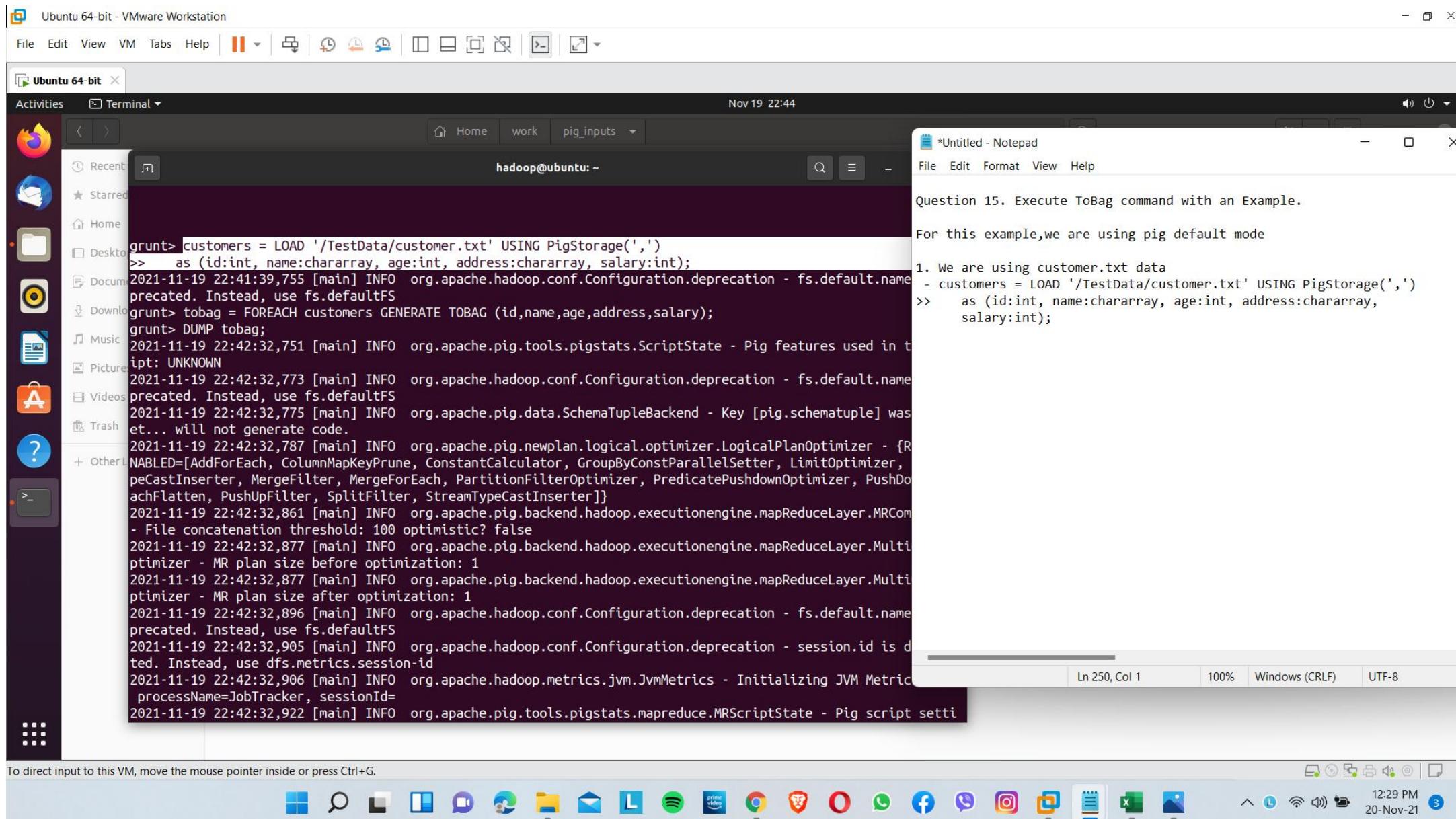
To direct input to this VM, move the mouse pointer inside or press Ctrl+G.

File Home Work pig_inputs

12:28 PM

20-Nov-21 3





File Edit View VM Help | || | | | | | | |

Ubuntu 64-bit

Activities Terminal

Nov 19 22:46

◻ ◻ ◻

Home work pig_inputs

Q E

```
hadoop@ubuntu: ~
grunt> customers = LOAD '/TestData/customer.txt' USING PigStorage(',');
-> as (id:int, name:chararray, age:int, address:chararray, salary:int);
2021-11-19 22:41:39,755 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name
precated. Instead, use fs.defaultFS
grunt> tobag = FOREACH customers GENERATE TOBAG (id,name,age,address,salary);
grunt> DUMP tobag;
2021-11-19 22:42:32,751 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in t
ip: UNKNOWN
2021-11-19 22:42:32,773 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name
precated. Instead, use fs.defaultFS
2021-11-19 22:42:32,775 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was
et... will not generate code.
2021-11-19 22:42:32,787 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - {R
UNABLE=[AddForEach, ColumnMapKeyPrune, ConstantCalculator, GroupByConstParallelSetter, LimitOptimizer,
peCastInserter, MergeFilter, MergeForEach, PartitionFilterOptimizer, PredicatePushdownOptimizer, PushDo
achFlatten, PushUpFilter, SplitFilter, StreamTypeCastInserter]}
2021-11-19 22:42:32,861 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCom
- File concatenation threshold: 100 optimistic? false
2021-11-19 22:42:32,877 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.Multi
ptimizer - MR plan size before optimization: 1
2021-11-19 22:42:32,877 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.Multi
ptimizer - MR plan size after optimization: 1
2021-11-19 22:42:32,896 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name
precated. Instead, use fs.defaultFS
2021-11-19 22:42:32,905 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - session.id is d
ted. Instead, use dfs.metrics.session-id
2021-11-19 22:42:32,906 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Initializing JVM Metrics
processName=JobTracker, sessionId=
2021-11-19 22:42:32,922 [main] INFO org.apache.pig.tools.pigstats.mapreduce.MRScriptState - Pig script setti
```

*Untitled - Notepad

File Edit Format View Help

Question 15. Execute ToBag command with an Example.

For this example, we are using pig default mode

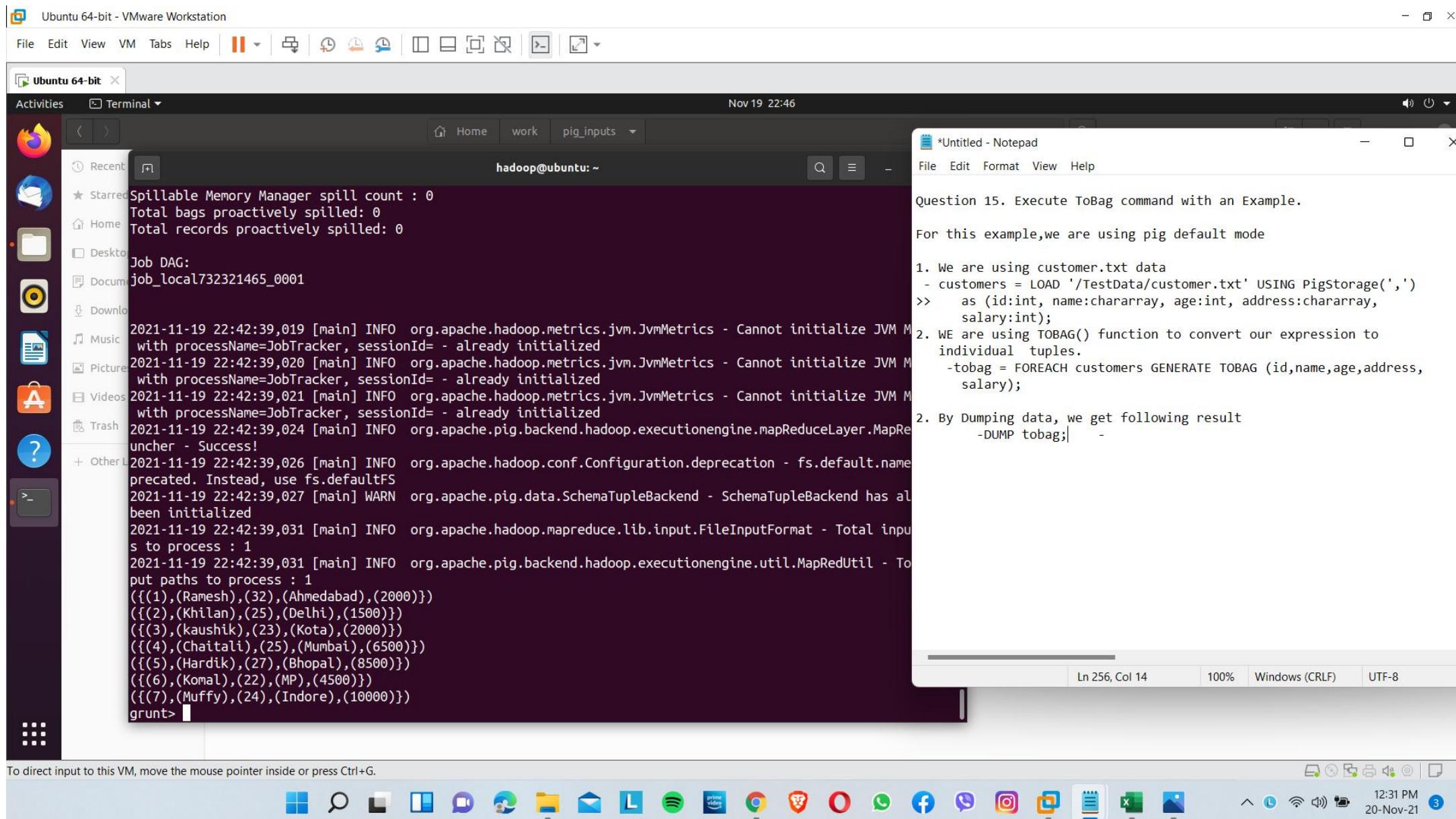
1. We are using customer.txt data
 - customers = LOAD '/TestData/customer.txt' USING PigStorage(',')
-> as (id:int, name:chararray, age:int, address:chararray, salary:int);
 2. WE are using TOBAG() function to convert our expression to individual tuples.
 - tobag = FOREACH customers GENERATE TOBAG (id,name,age,address, salary);
2. By Dumping data, we get following result
-DUMP tobag;

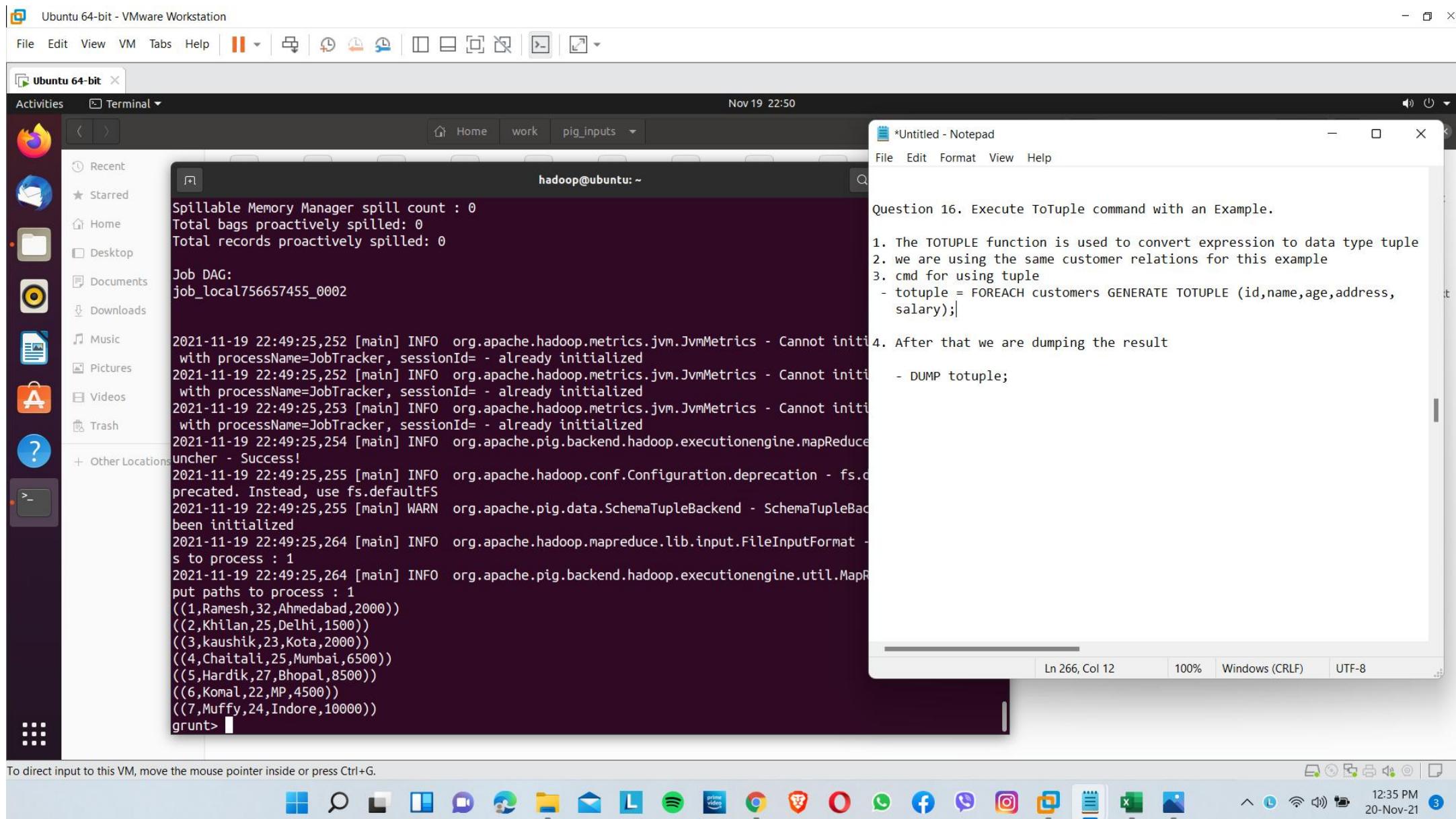
To direct input to this VM, move the mouse pointer inside or press Ctrl+G.



◻ ◻ ◻

12:31 PM
20-Nov-21 3





File Edit View VM Help | || | | | | | | | |

Ubuntu 64-bit

Activities Terminal

Nov 19 22:54



```
grunt> tomap = FOREACH customers GENERATE TOMAP(name, address);
grunt> DUMP tomap;
2021-11-19 22:53:26,037 [main] INFO  org.apache.pig.tools.pigstats.ScriptState - Pig feature
ipt: UNKNOWN
2021-11-19 22:53:26,045 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - fs.co
precated. Instead, use fs.defaultFS
2021-11-19 22:53:26,046 [main] INFO  org.apache.pig.data.SchemaTupleBackend - Key [pig.schem
et... will not generate code.
2021-11-19 22:53:26,046 [main] INFO  org.apache.pig.newplan.logical.optimizer.LogicalPlanOpti
NABLED=[AddForEach, ColumnMapKeyPrune, ConstantCalculator, GroupByConstParallelSetter, Limit
peCastInserter, MergeFilter, MergeForEach, PartitionFilterOptimizer, PredicatePushdownOptim
achFlatten, PushUpFilter, SplitFilter, StreamTypeCastInserter]
2021-11-19 22:53:26,048 [main] INFO  org.apache.pig.newplan.logical.rules.ColumnPruneVisitor
for customers: $0, $2, $4
2021-11-19 22:53:26,049 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduce
- File concatenation threshold: 100 optimistic? false
2021-11-19 22:53:26,050 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduce
ptimizer - MR plan size before optimization: 1
2021-11-19 22:53:26,050 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduce
ptimizer - MR plan size after optimization: 1
2021-11-19 22:53:26,054 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - fs.co
precated. Instead, use fs.defaultFS
2021-11-19 22:53:26,055 [main] INFO  org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initia
with processName=JobTracker, sessionId= - already initialized
2021-11-19 22:53:26,056 [main] INFO  org.apache.pig.tools.pigstats.mapreduce.MRScriptState -
ngs are added to the job
2021-11-19 22:53:26,056 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduce
ompiler - mapred.job.reduce.markreset.buffer.percent is not set, set to default 0.3
2021-11-19 22:53:26,056 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlC
ompiler - This job cannot be converted run in-process
2021-11-19 22:53:26,090 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlC
```

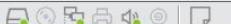
*Untitled - Notepad

File Edit Format View Help

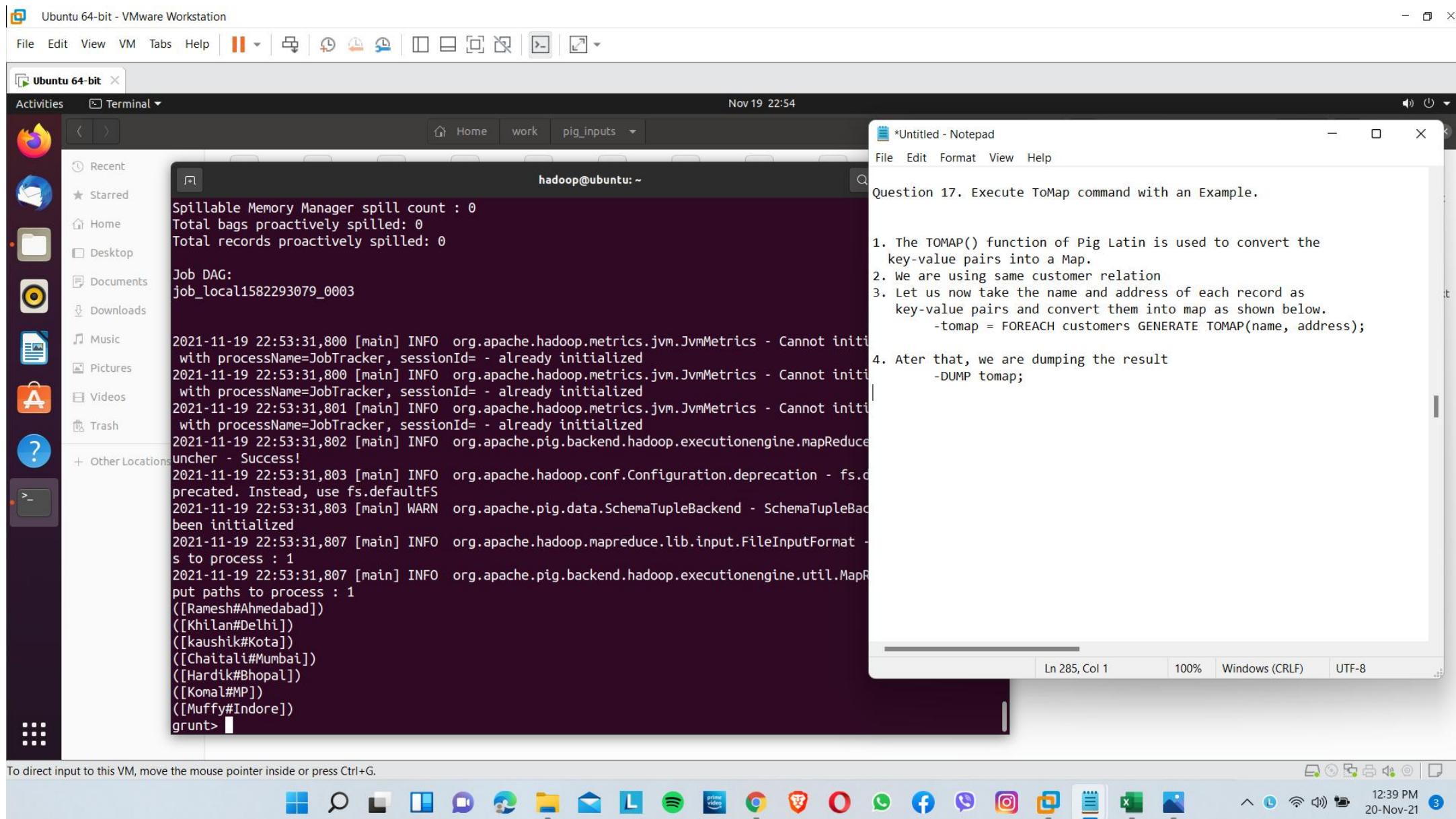
Question 17. Execute ToMap command with an Example.

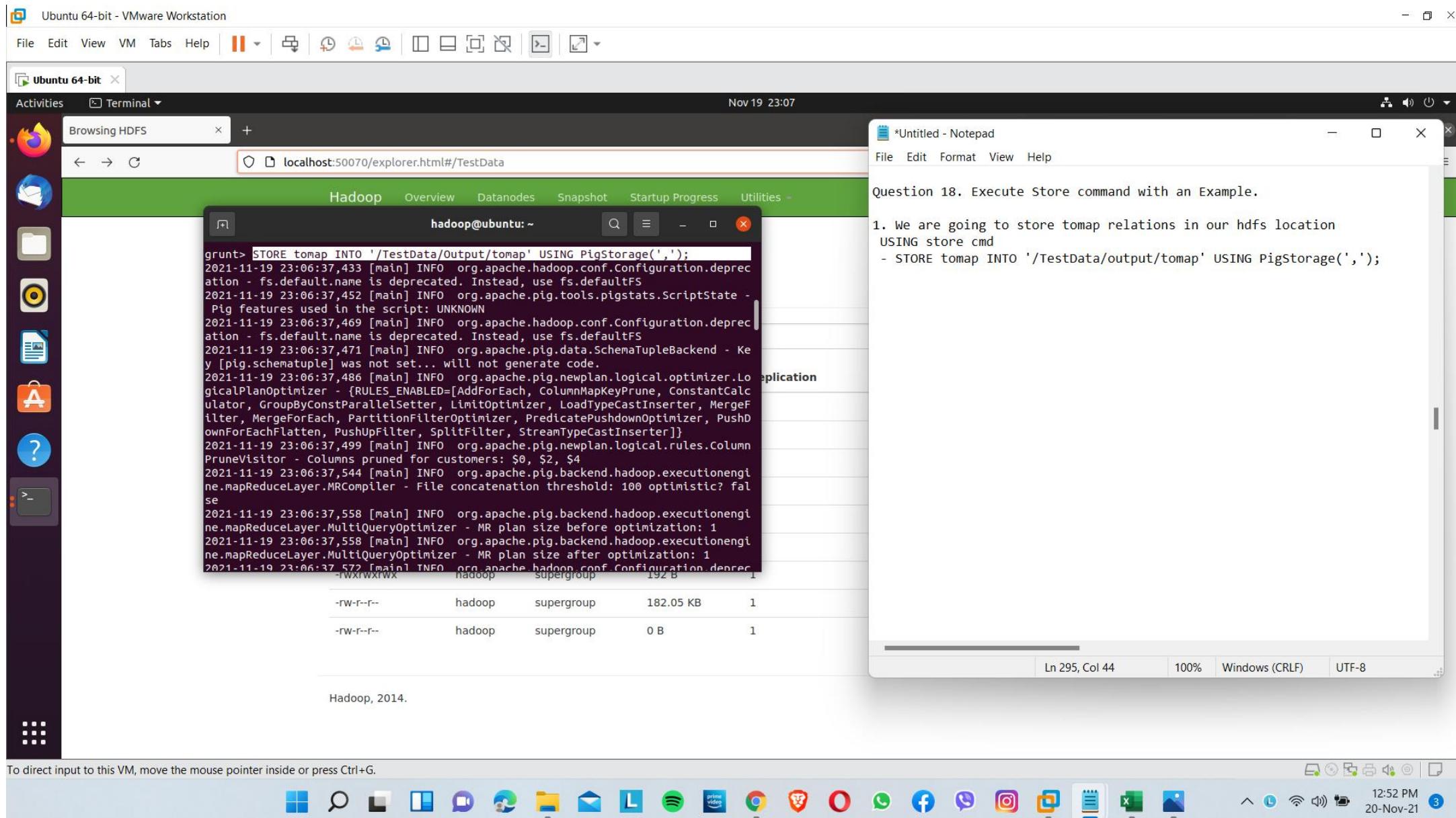
1. The TOMAP() function of Pig Latin is used to convert the key-value pairs into a Map.
2. We are using same customer relation
3. Let us now take the name and address of each record as key-value pairs and convert them into map as shown below.
-tomap = FOREACH customers GENERATE TOMAP(name, address);
4. After that, we are dumping the result
-DUMP tomap;

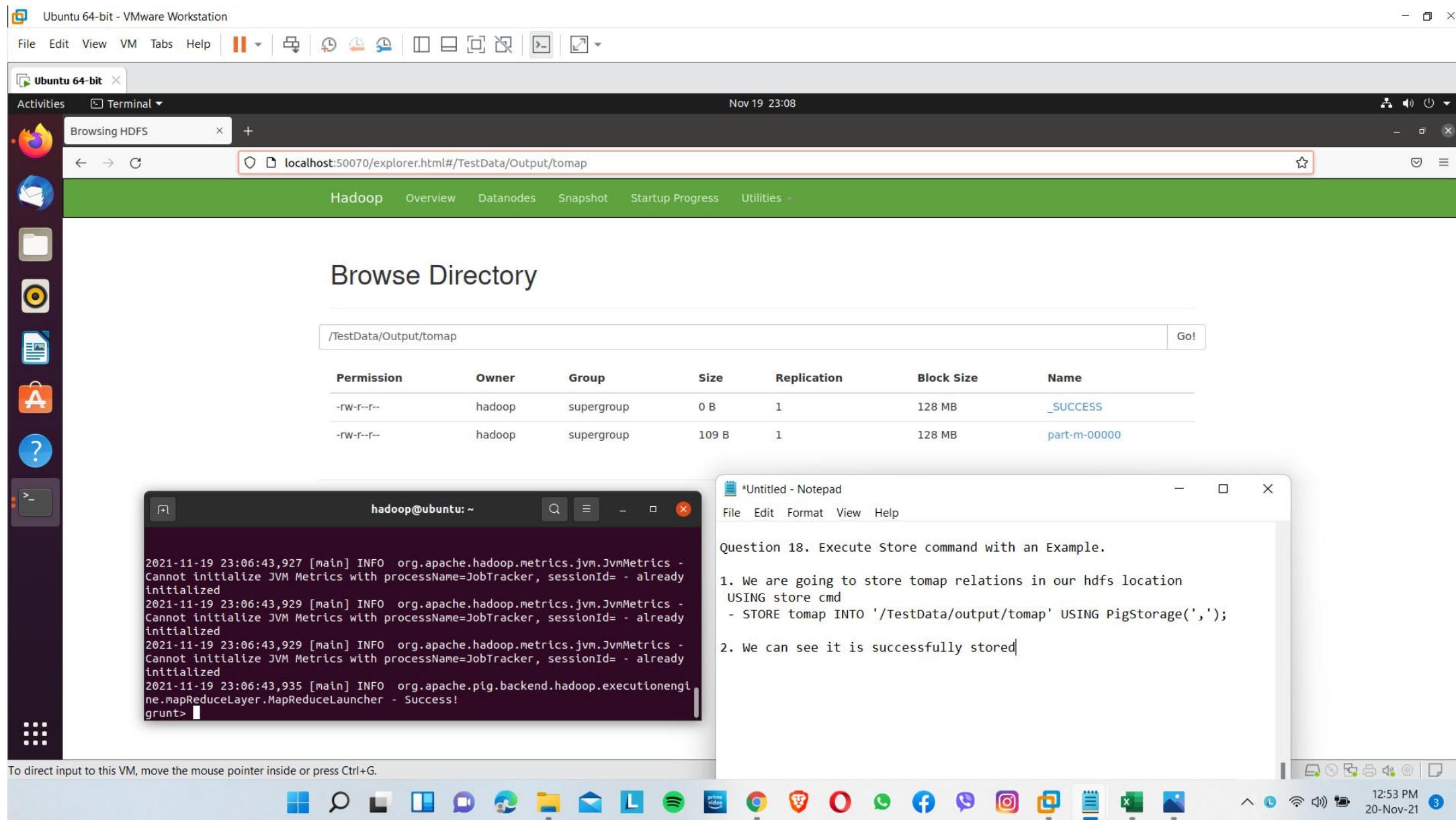
To direct input to this VM, move the mouse pointer inside or press Ctrl+G.

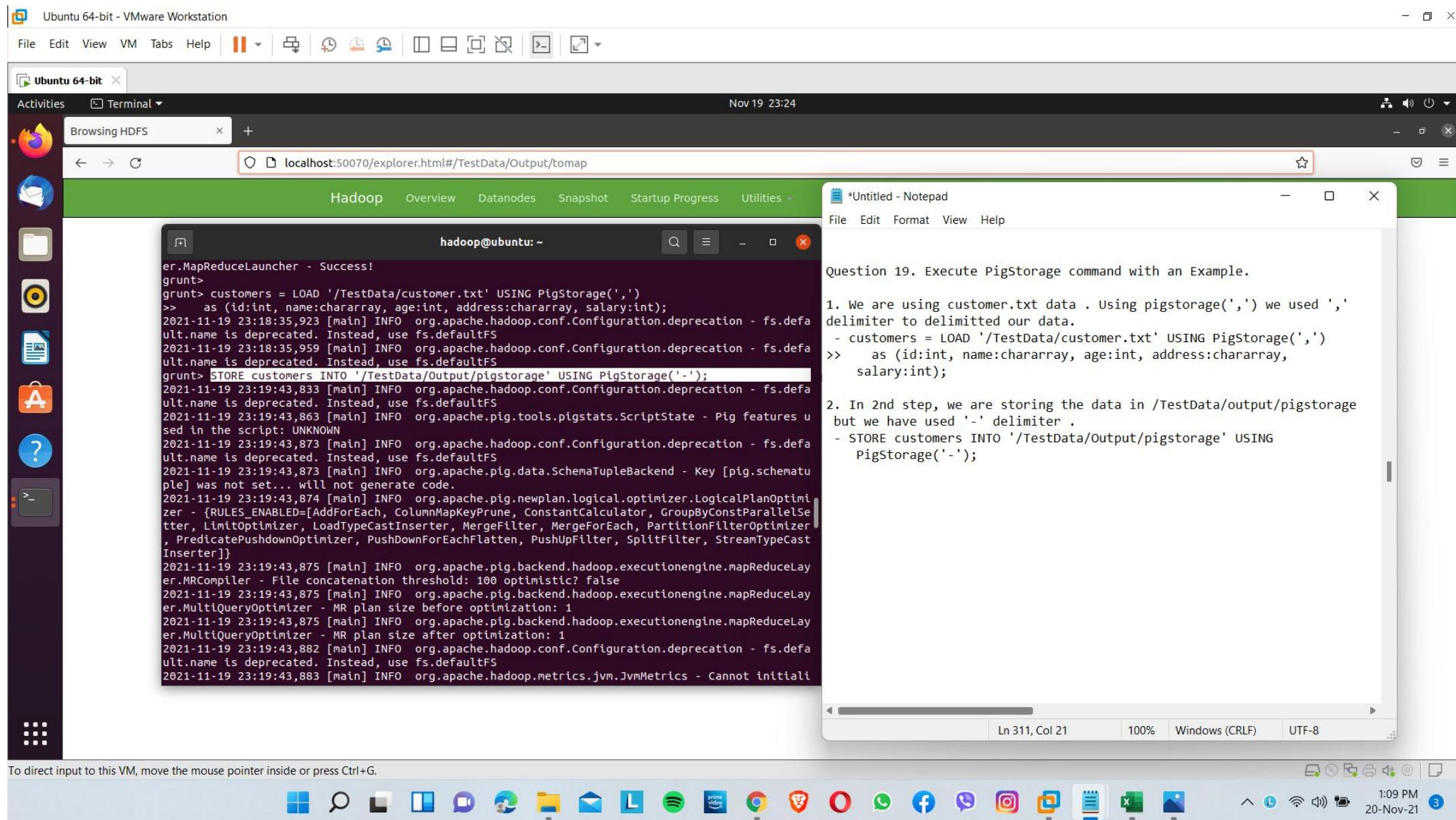


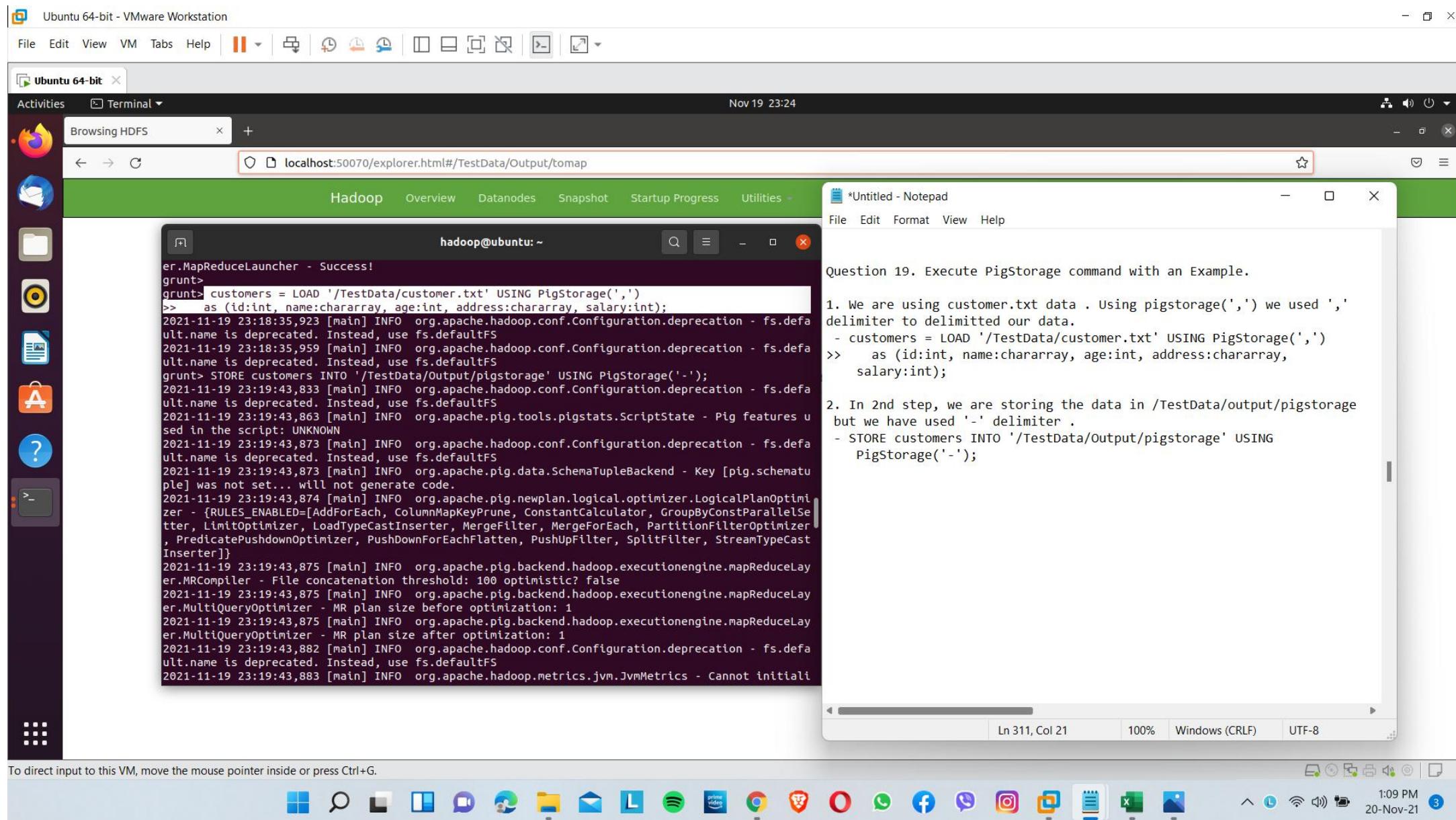
12:39 PM 20-Nov-21 3

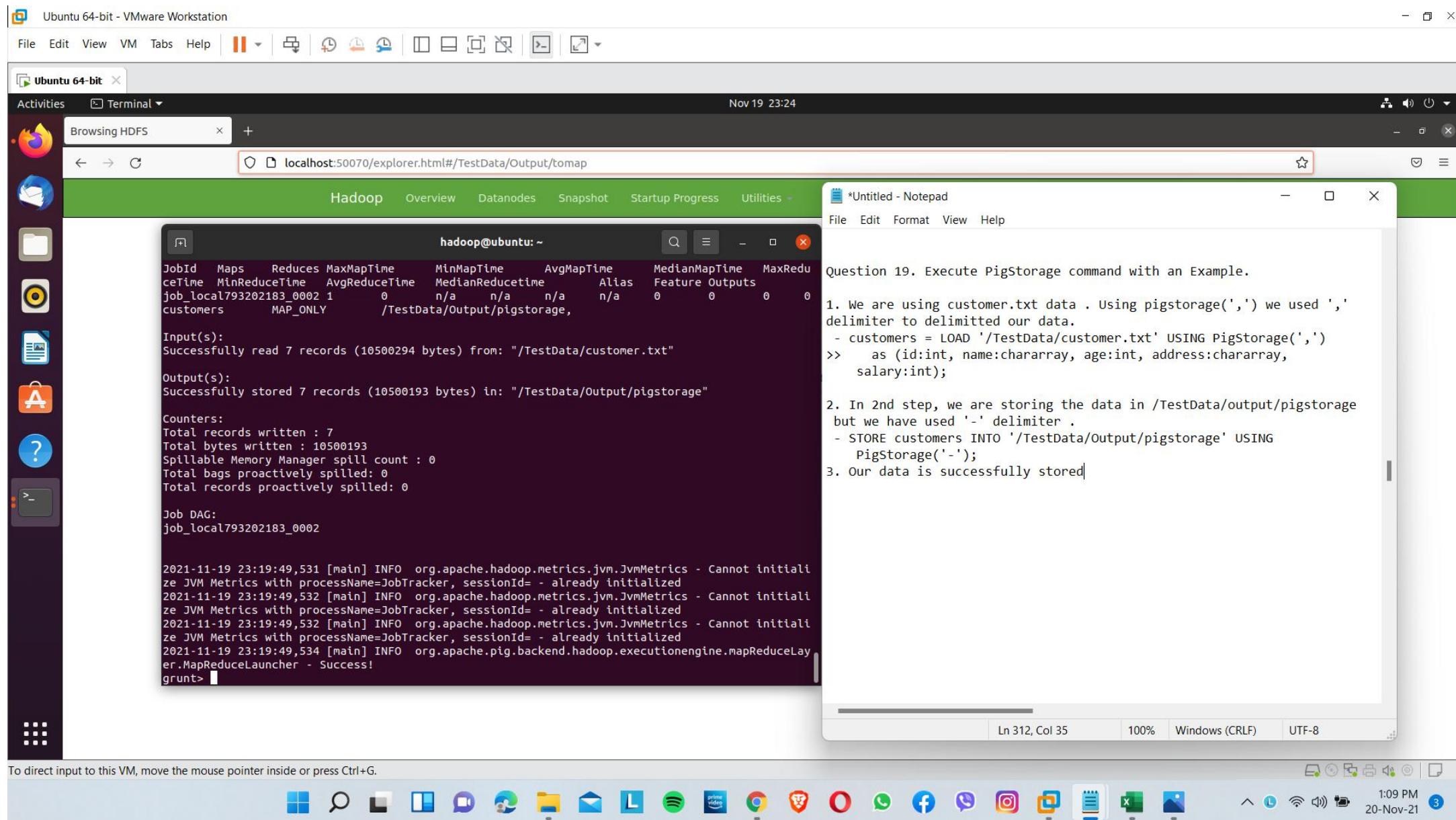












File Edit View VM Tabs Help | || ⌂ ⌃ ⌄ ⌅ ⌆ ⌇ ⌈ ⌉ ⌊ ⌋

Ubuntu 64-bit

Activities Text Editor

Nov 19 23:25

part-m-00000

~/Downloads

Save

☰

-

□

×

1 1-Ramesh-32-Ahmedabad-2000
2 2-Khilan-25-Delhi-1500
3 3-kaushik-23-Kota-2000
4 4-Chaitali-25-Mumbai-6500
5 5-Hardik-27-Bhopal-8500
6 6-Komal-22-MP-4500
7 7-Muffy-24-Indore-10000

*Untitled - Notepad

File Edit Format View Help

Question 19. Execute PigStorage command with an Example.

1. We are using customer.txt data . Using pigstorage(',') we used ',' delimiter to delimitted our data.

```
- customers = LOAD '/TestData/customer.txt' USING PigStorage(',')  
-> as (id:int, name:chararray, age:int, address:chararray,  
      salary:int);
```

2. In 2nd step, we are storing the data in /TestData/output/pigstorage but we have used '-' delimiter .

```
- STORE customers INTO '/TestData/output/pigstorage' USING  
  PigStorage('-');
```

3. Our data is successfully stored

4. This is our output data

Ln 313, Col 27

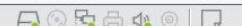
100%

Windows (CRLF)

UTF-8

Plain Text Tab Width: 8 Ln 1, Col 1 INS

To direct input to this VM, move the mouse pointer inside or press Ctrl+G.



1:10 PM

20-Nov-21

3

```
1 1-Ramesh-32-Ahmedabad-2000
2 2-Khilan-25-Delhi-1500
3 3-kaushik-23-Kota-2000
4 4-Chaitali-25-Mumbai-6500
```

hadoop@ubuntu: ~

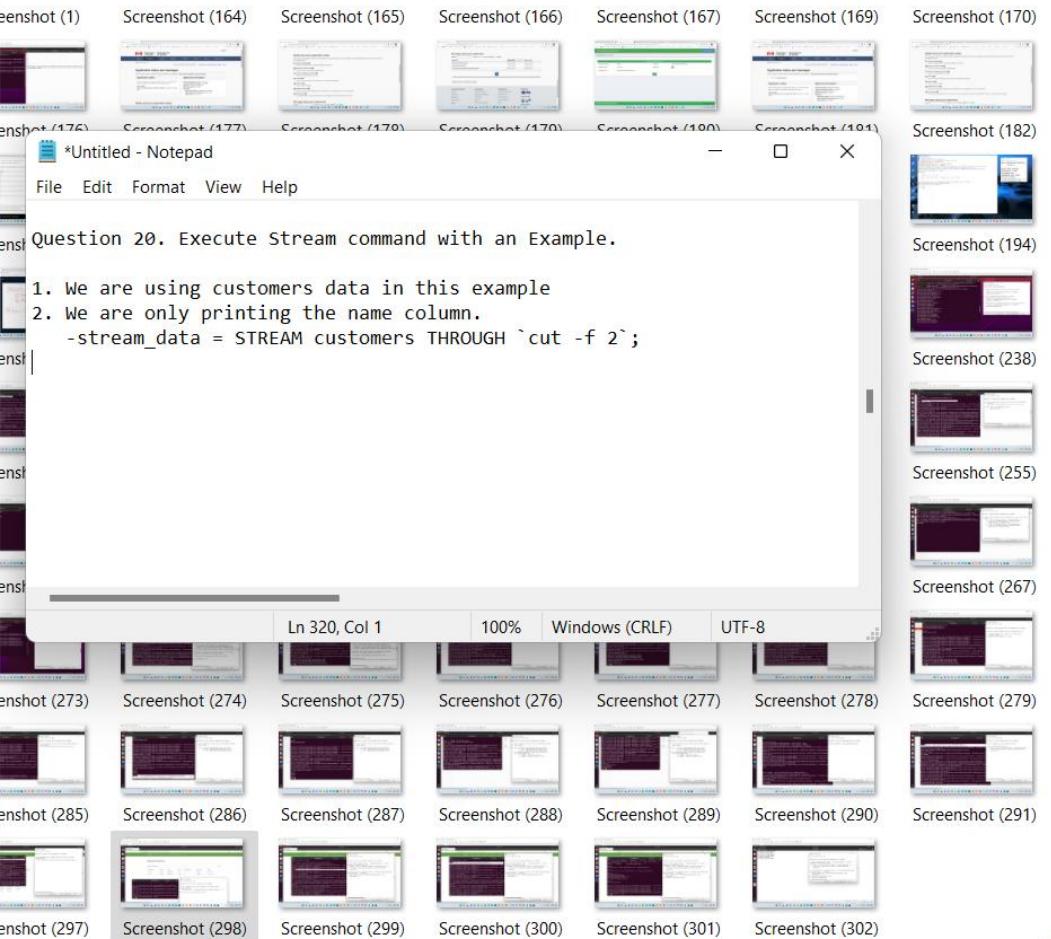
```
grunt> stream data = STREAM customers THROUGH `cut -f 2`;
grunt> DUMP stream_data;
2021-11-19 23:31:13,621 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: STREAMING
2021-11-19 23:31:13,644 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-11-19 23:31:13,647 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schema] was not set... will not generate code.
2021-11-19 23:31:13,662 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - {RULES_ENABLED=[AddForEach, ColumnMapKeyPrune, ConstantCalculator, GroupByConstParalleliser, LimitOptimizer, LoadTypeCastInserter, MergeFilter, MergeForEach, PartitionFilterOptimizer, PredicatePushdownOptimizer, PushDownForEachFlatten, PushUpFilter, SplitFilter, StreamTypeCastInserter]}
2021-11-19 23:31:13,719 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCompiler - File concatenation threshold: 100 optimistic? false
2021-11-19 23:31:13,730 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size before optimization: 1
2021-11-19 23:31:13,730 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size after optimization: 1
2021-11-19 23:31:13,743 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-11-19 23:31:13,749 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - session.id is deprecated. Instead, use dfs.metrics.session-id
2021-11-19 23:31:13,750 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Initializing JVM Metrics with processName=JobTracker, sessionId=
2021-11-19 23:31:13,762 [main] INFO org.apache.pig.tools.pigstats.mapreduce.MRScriptState - Pig script settings are added to the job
2021-11-19 23:31:13,765 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.reduce.markreset.buffer.percent is deprecated. Instead, use mapreduce.reduce.markreset.buffer.percent
```

To direct input to this VM, move the mouse pointer inside or press Ctrl+G.

Rotate left

Rotate right

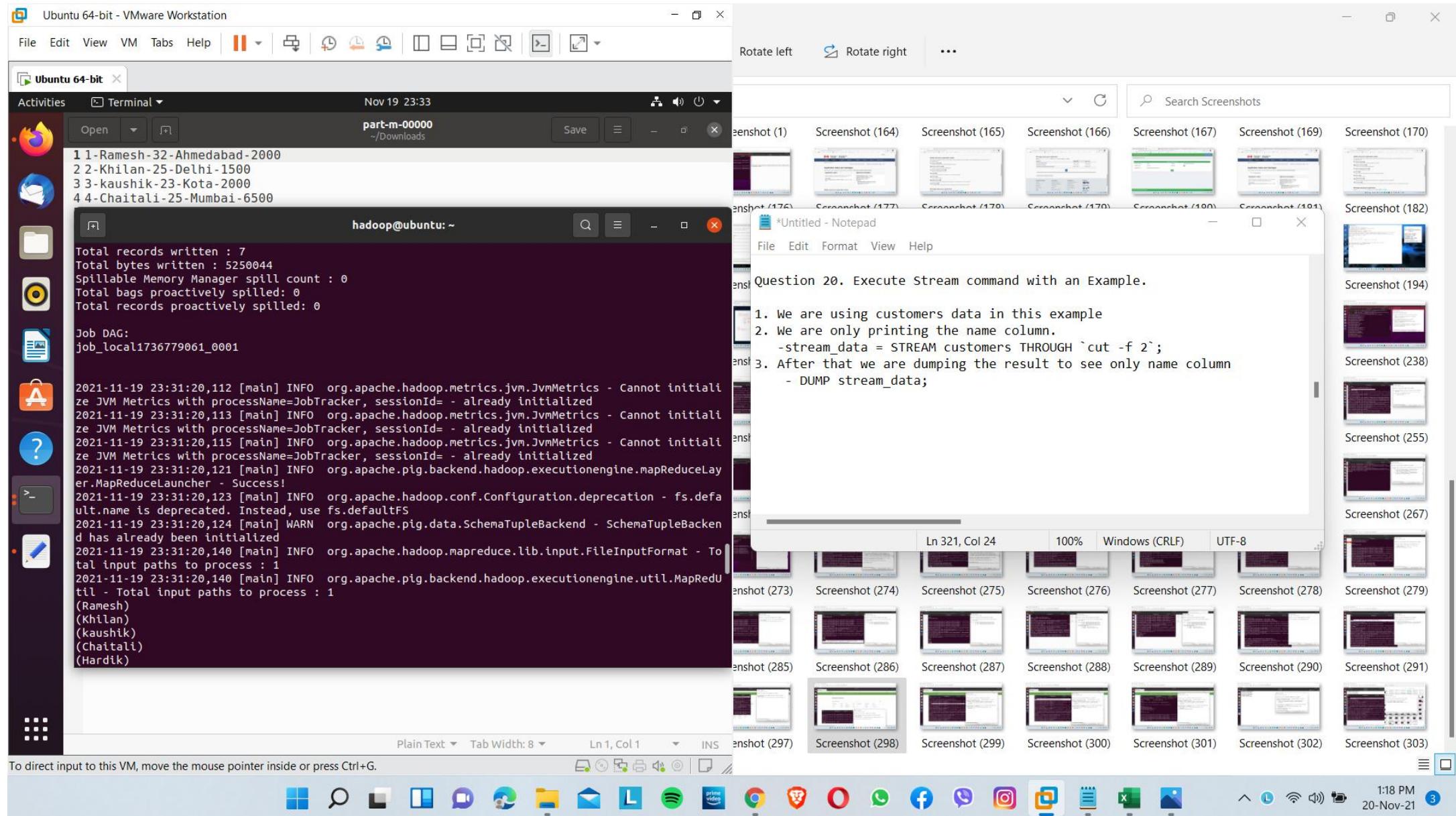
...



Question 20. Execute Stream command with an Example.

1. We are using customers data in this example
2. We are only printing the name column.

```
-stream_data = STREAM customers THROUGH `cut -f 2`;
```



File Edit View VM Tabs Help

Ubuntu 64-bit

Activities Terminal

Nov 19 23:34

part-m-00000

~/Downloads

Save

☰

-

□

X

Rotate left

Rotate right

...

Screenshot (1)

Screenshot (164)

Screenshot (165)

Screenshot (166)

Screenshot (167)

Screenshot (169)

Screenshot (170)

Screenshot (176)

Screenshot (177)

Screenshot (178)

Screenshot (179)

Screenshot (180)

Screenshot (181)

Screenshot (182)

Screenshot (183)

Screenshot (184)

Screenshot (185)

Screenshot (186)

Screenshot (187)

Screenshot (188)

Screenshot (189)

Screenshot (190)

Screenshot (191)

Screenshot (192)

Screenshot (193)

Screenshot (194)

Screenshot (195)

Screenshot (196)

Screenshot (197)

Screenshot (198)

Screenshot (199)

Screenshot (200)

Screenshot (201)

Screenshot (202)

Screenshot (203)

Screenshot (204)

Screenshot (205)

Screenshot (206)

Screenshot (207)

Screenshot (208)

Screenshot (209)

Screenshot (210)

Screenshot (211)

Screenshot (212)

Screenshot (213)

Screenshot (214)

Screenshot (215)

Screenshot (216)

Screenshot (217)

Screenshot (218)

Screenshot (219)

Screenshot (220)

Screenshot (221)

Screenshot (222)

Screenshot (223)

Screenshot (224)

Screenshot (225)

Screenshot (226)

Screenshot (227)

Screenshot (228)

Screenshot (229)

Screenshot (230)

Screenshot (231)

Screenshot (232)

Screenshot (233)

Screenshot (234)

Screenshot (235)

Screenshot (236)

Screenshot (237)

Screenshot (238)

Screenshot (239)

Screenshot (240)

Screenshot (241)

Screenshot (242)

Screenshot (243)

Screenshot (244)

Screenshot (245)

Screenshot (246)

Screenshot (247)

Screenshot (248)

Screenshot (249)

Screenshot (250)

Screenshot (251)

Screenshot (252)

Screenshot (253)

Screenshot (254)

Screenshot (255)

Screenshot (256)

Screenshot (257)

Screenshot (258)

Screenshot (259)

Screenshot (260)

Screenshot (261)

Screenshot (262)

Screenshot (263)

Screenshot (264)

Screenshot (265)

Screenshot (266)

Screenshot (267)

Screenshot (268)

Screenshot (269)

Screenshot (270)

Screenshot (271)

Screenshot (272)

Screenshot (273)

Screenshot (274)

Screenshot (275)

Screenshot (276)

Screenshot (277)

Screenshot (278)

Screenshot (279)

Screenshot (280)

Screenshot (281)

Screenshot (282)

Screenshot (283)

Screenshot (284)

Screenshot (285)

Screenshot (286)

Screenshot (287)

Screenshot (288)

Screenshot (289)

Screenshot (290)

Screenshot (291)

Screenshot (292)

Screenshot (293)

Screenshot (294)

Screenshot (295)

Screenshot (296)

Screenshot (297)

Screenshot (298)

Screenshot (299)

Screenshot (300)

Screenshot (301)

Screenshot (302)

Screenshot (303)

Screenshot (304)

Screenshot (305)

Screenshot (306)

Screenshot (307)

Screenshot (308)

Screenshot (309)

Screenshot (310)

Screenshot (311)

Screenshot (312)

Screenshot (313)

Screenshot (314)

Screenshot (315)

Screenshot (316)

Screenshot (317)

Screenshot (318)

Screenshot (319)

Screenshot (320)

Screenshot (321)

Screenshot (322)

Screenshot (323)

Screenshot (324)

Screenshot (325)

Screenshot (326)

Screenshot (327)

Screenshot (328)

Screenshot (329)

Screenshot (330)

Screenshot (331)

Screenshot (332)

Screenshot (333)

Screenshot (334)

Screenshot (335)

Screenshot (336)

Screenshot (337)

Screenshot (338)

Screenshot (339)

Screenshot (340)

Screenshot (341)

Screenshot (342)

Screenshot (343)

Screenshot (344)

Screenshot (345)

Screenshot (346)

Screenshot (347)

Screenshot (348)

Screenshot (349)

Screenshot (350)

To direct input to this VM, move the mouse pointer inside or press Ctrl+G.


 1:19 PM
 20-Nov-21
 3

```
1 1-Ramesh-32-Ahmedabad-2000  
2 2-Khilan-25-Delhi-1500  
3 3-kaushik-23-Kota-2000  
4 4-Chaitali-25-Mumbai-6500
```

hadoop@ubuntu: ~

```
Total bags proactively spilled: 0  
Total records proactively spilled: 0
```

```
Job DAG:  
job_local1911418541_0002
```

```
2021-11-19 23:31:40,802 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize JVM Metrics with processName=JobTracker, sessionId= - already initialized  
2021-11-19 23:31:40,803 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize JVM Metrics with processName=JobTracker, sessionId= - already initialized  
2021-11-19 23:31:40,803 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize JVM Metrics with processName=JobTracker, sessionId= - already initialized  
2021-11-19 23:31:40,805 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!  
2021-11-19 23:31:40,805 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS  
2021-11-19 23:31:40,805 [main] WARN org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend has already been initialized  
2021-11-19 23:31:40,810 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1  
2021-11-19 23:31:40,810 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapReduceUtil - Total input paths to process : 1  
(Ahmedabad)  
(Delhi)  
(Kota)  
(Mumbai)  
(Bhopal)  
(MP)  
(Indore)  
grunt> |
```

Plain Text

Tab Width: 8

Ln 1, Col 1

INS

-

Screenshot (297)

Screenshot (298)

Screenshot (299)

Screenshot (300)

Screenshot (301)

Screenshot (302)

Screenshot (303)

To direct input to this VM, move the mouse pointer inside or press Ctrl+G.

