

Issued on: March 05, 2016

Due by 11:30PM EST, March 11, 2016

Submission is late by 40 minutes

Please, describe every step of your work and present all intermediate and final results in a Word document. Please, copy past text version of all essential command and snippets of results into the Word document. We cannot retype text that is in JPG images. Please, always submit a copy of original, working scripts and/or class files you used as separate files. Sometimes we need to run your code and retyping is too costly. Please, submit to the class drop box. For issues and comments visit the class Discussion Board. You can solve these problems using any language of your choice.

Problem 1. Go to an online newspaper and select paragraphs from two articles in a similar field, about politics, art, movies, or any other topic of your choice. Let those paragraphs be moderately small, a few lines and around 100 words.

Save those paragraphs as .txt files and then import them into two Spark RDD objects, paragraphA and paragraphB.

Use Spark transformation functions to transform those initial RDD-s into RDD-s that contain only words.

Solution:

I import and split the files into words using a class mySplitter, which is saved as a stand alone file named mySplitter.java

```
class mySplitter implements Function<JavaRDD<String>,JavaRDD<String> >
{
    private String sregex;

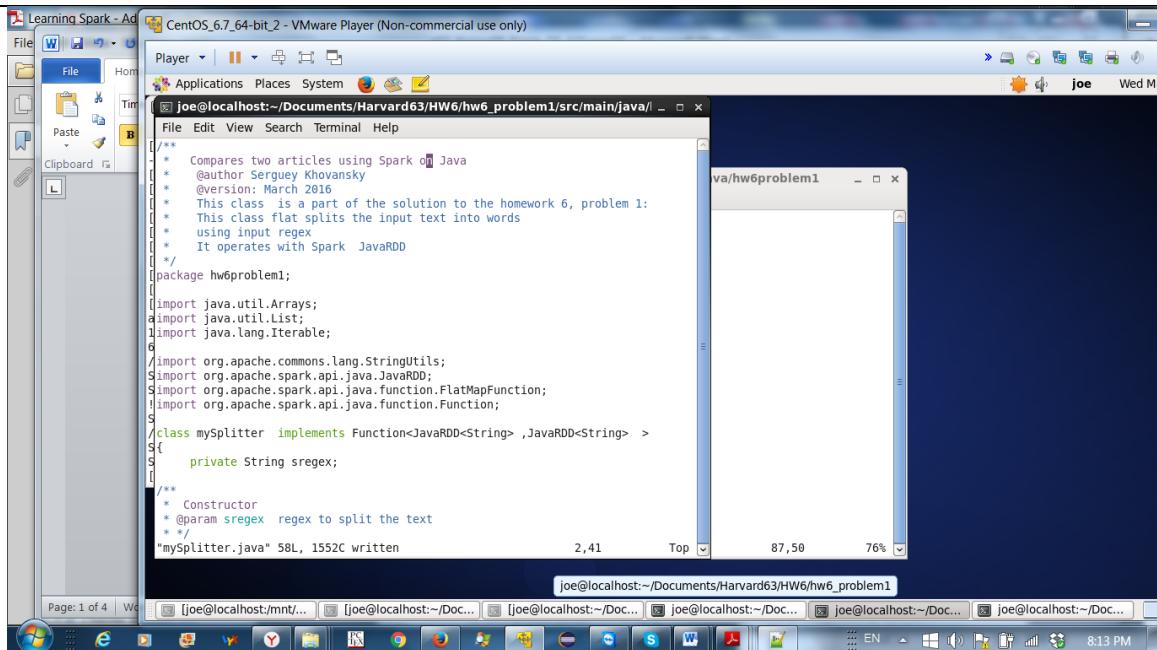
    /**
     * Constructor
     * @param sregex regex to split the text
     */
    mySplitter(String sregex)
    {
        this.sregex = sregex;
    }

    /**
     * Call function is implemented from Function class of Spark
     * @param input JavaRDD<String> object, obtained from input text file
     * @return words JavaRDD<String> output object with words split
     */
    @Override
    public JavaRDD<String> call( JavaRDD<String> input) {
        JavaRDD<String> words = input.flatMap
        (
            new FlatMapFunction<String, String>()
```

```

    {
        public Iterable<String> call(String x)
        {
            //return Arrays.asList(x.toLowerCase().split( "\\\W+" ) );
            return Arrays.asList(x.toLowerCase().split(sregex));
        }
    );
    return words;
}

```



To read and split the files use the following lines in the class with main method (named hw6problem1)

```

String sregex = "\\\W+";
MySplitter MySplitter = new MySplitter(sregex );

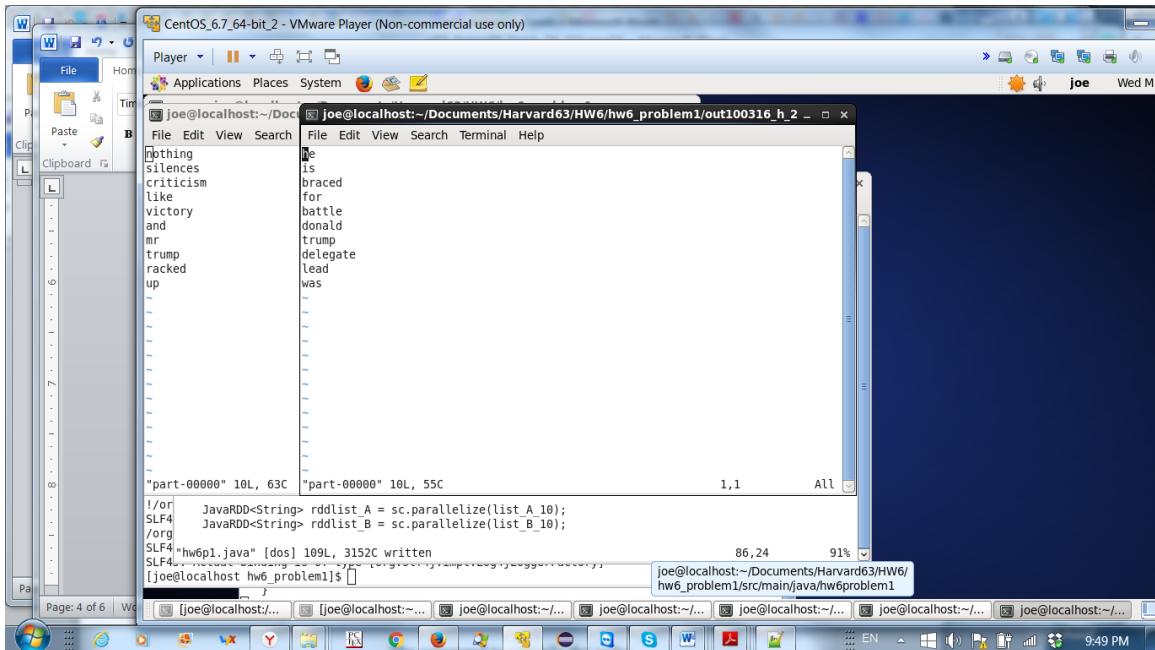
JavaRDD<String> words1 = MySplitter.call(input1);
JavaRDD<String> words2 = MySplitter.call(input2);

```

List for us the first 10 words in each RDD.

It is possible to create a special class for each operation, similar to mySplitter class, which I used to split the original input files, but I think that the overhead in that class would exceed the several extra lines that I need to print out here to handle the output twice.

```
List<String> list_A_10 = paragraphA.take(10);
List<String> list_B_10 = paragraphB.take(10);
JavaRDD<String> rddlist_A = sc.parallelize(list_A_10);
JavaRDD<String> rddlist_B = sc.parallelize(list_B_10);
rddlist_A.saveAsTextFile(outputFile1);
rddlist_B.saveAsTextFile(outputFile2);
```



Subsequently create RDD-s that contain only unique words in each of paragraphs.

Relevant line:

```
JavaRDD<String> rddUniqueA_B
=paragraphB.subtract(paragraphA).union(paragraphA.subtract(paragraphB)).distinct();
```

Output:

```
decisive
obama
silences
themselves
week
opponents
cruz
prize
one
attack
```

```
certain  
big  
national  
ads  
several  
only  
been  
poised  
still  
first  
losing  
she  
whites
```

```
joe@localhost:~/Documents/Harvard63/HW6/hw6_problem1/src/main/java/hw6problem1
```

```
List<String> list_B_10 = paragraphB.take(10);  
  
JavaRDD<String> rddList_A = sc.parallelize(list_A_10);  
JavaRDD<String> rddList_B = sc.parallelize(list_B_10);  
  
JavaRDD<String> rddUniqueA_B = paragraphB.subtract(paragraphA).union(paragraphA.subtract(paragraphB)).distinct();  
  
JavaRDD<String> rddUniqueAonly = paragraphA.distinct().subtract(paragraphB).distinct();  
JavaRDD<String> rddIntersection_AB = paragraphA.intersection(paragraphB);  
  
String outputFile_union = outputFile+"u";  
String outputFile_onlyA = outputFile+"only";  
String outputFile_ABintersec = outputFile+"ABintersec";  
  
rddList_A.saveAsTextFile(outputFile1);  
rddList_B.saveAsTextFile(outputFile2);  
rddUniqueA_B.saveAsTextFile(outputFile_union);  
rddUniqueAonly.saveAsTextFile(outputFile_onlyA);  
rddIntersection_AB.saveAsTextFile(outputFile_ABintersec);  
  
"hw6p1.java" [dos] 125L, 3792C written
```

```
104,1 94%
```

```
[joe@localhost hw6_problem1]$
```

```
joe@localhost:~/Documents/Harvard63/HW6/hw6_problem1/src/main/java/hw6problem1
```

Then create an RDD that contains only words that are present in paragraphA but not in paragraphB.

Relevant part of the code:

```
JavaRDD<String> rddUniqueAonly = paragraphA.distinct().subtract(paragraphB).distinct();
```

Picture of the relevant code with the part of the output


```

campaign
he
michigan
the
his
tuesday
was
won
but
mississippi
is
trump
had
on
with
in
from
well
that
a
to
of
and
"part-00000" 23L, 112C
"/org/slf4j/impl/StaticLoggerBinder.class"
SLF4J: Found binding in [jar:file:/usr/lib/flume-ng/lib/slf4j-log4j12-1.7.5.jar!
/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]
[joe@localhost hw6_problem1]$ 
```

Output

```

campaign
he
michigan
the
his
tuesday
was
won
but
mississippi
is
trump
had
on
with
in
from
well
that
a
to
of
and

```

Compilation:

Prepare Maven

```

<project>
  <groupId>hw6problem1</groupId>
  <artifactId>spark-example</artifactId>
  <modelVersion>4.0.0</modelVersion>
  <name>problem1</name>
  <packaging>jar</packaging>
  <version>0.0.2</version>
  <dependencies>
    <dependency><!-- Spark dependency --&gt;
      &lt;groupId&gt;org.apache.spark&lt;/groupId&gt;
      &lt;artifactId&gt;spark-core_2.10&lt;/artifactId&gt;
      &lt;version&gt;1.1.0&lt;/version&gt;
      &lt;scope&gt;provided&lt;/scope&gt;
    &lt;/dependency&gt;
  &lt;/dependencies&gt;
  &lt;properties&gt;
    &lt;java.version&gt;1.8&lt;/java.version&gt;
  &lt;/properties&gt;
  &lt;build&gt;
    &lt;pluginManagement&gt;</pre>

```

```

<plugins>
  <plugin>
    <groupId>org.apache.maven.plugins</groupId>
    <artifactId>maven-compiler-plugin</artifactId>
    <version>3.1</version>
    <configuration>
      <source>${java.version}</source>
      <target>${java.version}</target>
    </configuration>
  </plugin>
</plugins>
</pluginManagement>
</build>
</project>

```

Place **pom.xml** in the folder :

```
/home/joe/Documents/Harvard63/HW6/hw6_problem1
[joe@localhost hw6_problem1]$ ls
hw6p1_data pom.xml src target
```

Place the code **hw6p1.java** into the folder:

```
/home/joe/Documents/Harvard63/HW6/hw6_problem1/src/hw6problem1
[joe@localhost hw6problem1]$ ls
hw6p1.java
```

Run maven, being in the folder where **pom.xml** is located, i.e.

```
/home/joe/Documents/Harvard63/HW6/hw6_problem1
[joe@localhost hw6_problem1]$ mvn clean && mvn compile && mvn package
```

```

CentOS_6.7_64-bit_2 - VMware Player (Non-commercial use only)

joe@localhost:~/Documents/Harvard63/HW6/hw6_problem1
File Edit View Search Terminal Help
[INFO] Scanning for projects...
[INFO]
[INFO] -----
[INFO] Building problem1 0.0.2
[INFO] -----
[INFO] --- maven-clean-plugin:2.5:clean (default-clean) @ spark-example ---
[INFO] Deleting /home/joe/Documents/Harvard63/HW6/hw6_problem1/target
[INFO]
[INFO] BUILD SUCCESS
[INFO] -----
[INFO] Total time: 1.933 s
[INFO] Finished at: 2016-03-09T17:05:42-08:00
[INFO] Final Memory: 5M/29M
[INFO] -----
[INFO] Scanning for projects...
[INFO] -----
[INFO] Building problem1 0.0.2
[INFO] -----
[INFO] --- maven-resources-plugin:2.6:resources (default-resources) @ spark-example ---
[WARNING] Using platform encoding (UTF-8 actually) to copy filtered resources, i.e. build is platform dependent!
[INFO] skip non existing resourceDirectory /home/joe/Documents/Harvard63/HW6/hw6_problem1/src/main/resources
[INFO]
[INFO] --- maven-compiler-plugin:3.1:compile (default-compile) @ spark-example ...

```

```

joe@localhost:~/Documents/Harvard63/HW6/hw6_problem1/src
File Edit View Search Terminal Help
JavaRDD<String> words2 = MySplitter.call(input2);

/**
 * Transform into word and count and sort by Key, using method
 */

/*JavaPairRDD<String, Integer> counts = words.mapToPair(
  new PairFunction<String, String, Integer>(){
    public Tuple2<String, Integer> call(String x){
      return new Tuple2(x, 1);
    }
  }).reduceByKey(new Function2<Integer, Integer, Integer>(){
    public Integer call(Integer x, Integer y){ return x + y
  }
}

/*
 * Save the word count back out to a text file, causing evaluation
 */
words1.saveAsTextFile(outputFile1);
words2.saveAsTextFile(outputFile2);

```

Run:

Folder:

```
/home/joe/Documents/Harvard63/HW6/hw6_problem1
```

```
[joe@localhost hw6_problem1]$ spark-submit --class hw6problem1.hw6p1
./target/spark-example-0.0.2.jar
file:///home/joe/Documents/Harvard63/HW6/hw6_problem1/hw6p1_data/hw6p1_da
ta_a.txt
file:///home/joe/Documents/Harvard63/HW6/hw6_problem1/hw6p1_data/hw6p1_da
ta_b.txt file:///home/joe/Documents/Harvard63/HW6/hw6_problem1/out100316_k
```

Problem 2.

Consider attached file `emps.txt`. It contains: name, age and salary of three employees.

Create RDD `emps` by importing that file into Spark.

```
public static void main(String[] args) throws Exception {
    // read the file name from command line
    String inputFile = args[0];
    String outputFile = args[1];
    //Create a Java Spark Context sc
    SparkConf conf = new SparkConf().setAppName("hw6p2");
    JavaSparkContext sc = new JavaSparkContext(conf);
    // Create RDD emps
    JavaRDD<String> emps = sc.textFile(inputFile);
```

Next create a new RDD `emps_fields` by transforming the content of every line in RDD `emps` into a tuple with three individual elements by splitting the lines on commas.

The relevant part of the code is.

```
JavaRDD<Tuple3<String, String, String>> emps_f = emps.map
(
    new Function<String, Tuple3<String, String, String>>()
    {
        public Tuple3<String, String, String> call(String x)
        {
            return new Tuple3(x.split(",")[0], x.split(",")[1], x.split(",")[2]);
        }
    }
);
```

Now comes something new. Spark has a class `Row` and you need to import it in your script or program. `Row` comes from the same package as class `SQLContext`.

First add several lines to pom.xml to be able to "digest" `SQLContext`, that is located in the package `org.apache.spark.sql.SQLContext`

```

<dependency>
  <groupId>org.apache.spark</groupId>
  <artifactId>spark-sql_2.10 </artifactId>
  <version>1.5.1</version>
</dependency>

```

Add the following line in the code:

```
import org.apache.spark.sql.Row;
```

Row class creates rows with named and typed fields. You need to apply “constructor”

Row to every tuple in RDD emps_fields, like:

employees = emps_fields.map(lambda e: Row(name = e[0], age = int(e[1]), salary = float(e[2])))
e[0], e[1] and e[2] are the first, second and third elements of the tuple e
representing a row (line) in RDD emps_fields. Note that int and float are types of
fields in new rows.

```

JavaRDD<Row> employees = emps_f.map
(
    new Function< Tuple3<String, String, String> , Row>()
    {
        public Row call(Tuple3<String, String, String> x)
        {
            return RowFactory.create( x._1(), Integer.parseInt(x._2()), Float.parseFloat( x._3()) );
        }
    }
);

```

Newly create RDD employees is now made of Row elements and is ready to be transformed into a DataFrame. You generate a DataFrame by passing an RDD of Row elements to the method createDataFrame() of class SQLContext . Do it.

```

//import relevant packages
import org.apache.spark.sql.SQLContext;
import org.apache.spark.sql.types.*;
import org.apache.spark.sql.DataFrame;

// create schema of dataset

```

```

StructType schema = new StructType(new StructField [] {
    new StructField("name", DataTypes.StringType, false, Metadata.empty()),
    new StructField("age", DataTypes.IntegerType, false, Metadata.empty()),
    new StructField("salary", DataTypes.FloatType, false, Metadata.empty())
});
// create dataframe
SQLContext sqlContext = new SQLContext(sc);
DataFrame people = sqlContext.createDataFrame(employees, schema);

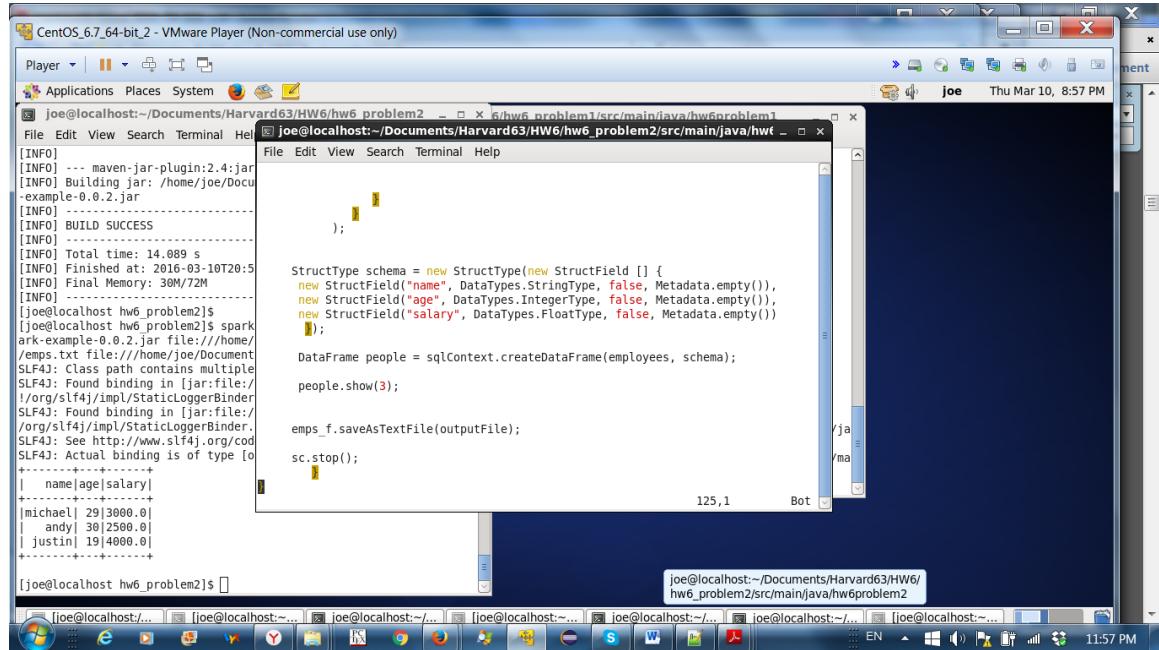
```

Show the content of new DataFrame.

```
people.show(3);
```

Output:

```
+-----+---+-----+
| name|age|salary|
+-----+---+-----+
|michael| 29|3000.0|
| andy| 30|2500.0|
| justin| 19|4000.0|
+-----+---+-----+
```

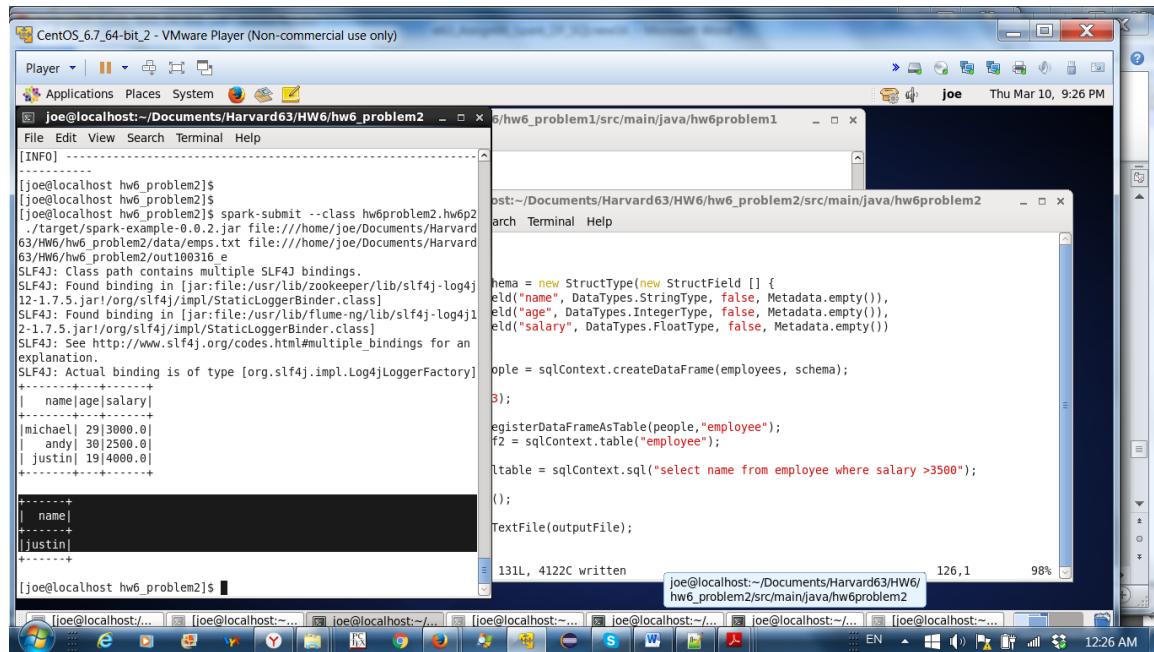


Transform this DataFrame into a Temporary Table and select names of all employees who have a salary greater than 3500.

```
sqlContext.registerDataFrameAsTable(people,"employee");
DataFrame df2 = sqlContext.table("employee");
DataFrame sqltable = sqlContext.sql("select name from employee where salary >3500");
sqltable.show();
```

Output:

```
+----+
| name|
+----+
|justin|
```



Running command

```
[joe@localhost hw6_problem2]$ spark-submit --class hw6problem2.hw6p2
./target/spark-example-0.0.2.jar
file:///home/joe/Documents/Harvard63/HW6/hw6_problem2/data/emps.txt
file:///home/joe/Documents/Harvard63/HW6/hw6_problem2/out100316_b
```

The code.

All detailed explanations are given in the discussion above

```

/**
 * @author Serguey Khovansky
 * @version: March 2016
 * This is solution of the homework 6, problem 2:
 *
 *
 * Consider attached file emps.txt. It contains: name, age and salary of three employees.
 * Create RDD emps by importing that file into Spark. Next create a new RDD emps_fields by
 * transforming the content of every line in RDD emps into a tuple with three individual
 * elements by splitting the lines on commas. Now comes something new. Spark has a class Row
 * and you need to import it in your script or program. Row comes from the same package as class
 * SQLContext. Row class creates rows with named and typed fields. You need to apply
 * "constructor" Row to every tuple in RDD emps_fields, like:
 * employees = emps_fields.map(lambda e: Row(name = e[0], age = int(e[1]), salary = float(e[2])))
 *
 * e[0], e[1] and e[2] are the first, second and third elements of the tuple e representing a
 * row (line) in RDD emps_fields. Note that int and float are types of fields in new rows.
 * Newly created RDD employees is now made of Row elements and is ready to be transformed into
 * a DataFrame. You generate a DataFrame by passing an RDD of Row elements to
 * the method createDataFrame() of class SQLContext. Do it. Show the content of new DataFrame.
 * Transform this DataFrame into a Temporary Table and select names of all
 * employees who have a salary greater than 3500.
 *
 */

package hw6problem2;

import java.util.Arrays;
import java.util.List;
import java.lang.Iterable;

import scala.Tuple3;

import org.apache.commons.lang.StringUtils;

import org.apache.spark.SparkConf;
import org.apache.spark.api.java.JavaRDD;
import org.apache.spark.api.java.JavaPairRDD;
import org.apache.spark.api.java.JavaSparkContext;
import org.apache.spark.api.java.function.FlatMapFunction;
import org.apache.spark.api.java.function.Function2;
import org.apache.spark.api.java.function.PairFunction;
import org.apache.spark.api.java.function.Function;

import org.apache.spark.sql.SQLContext;
import org.apache.spark.sql.Row;
import org.apache.spark.sql.RowFactory;

import org.apache.spark.sql.types.*;
import org.apache.spark.sql.DataFrame;

/**
 * The class hw6p1 make database operations with two text files using SPARK MapReduce technology
 */
public class hw6p2 {

    /**
     * @param args[] input and output paths and location of the file input and output files
     */
    public static void main(String[] args) throws Exception {
        String inputFile = args[0];
        String outputFile= args[1];

        /**
         * Create a Java Spark Context sc
         */
        SparkConf conf = new SparkConf().setAppName("hw6p2");
        JavaSparkContext sc = new JavaSparkContext(conf);

        // need for sql
    }
}

```

```

SQLContext sqlContext = new SQLContext(sc);

/**
 *      Create RDD emps
 */
JavaRDD<String> emps = sc.textFile(inputFile);

/**
 *Turn into Tuple of words
 */

JavaRDD<Tuple3<String, String, String>> emps_f = emps.map
(
    new Function<String, Tuple3<String, String, String>>()
    {
        public Tuple3<String, String, String> call(String x)
        {
            List L= Arrays.asList(x.toLowerCase().split( "\\\W+" ) );
            return new Tuple3(L.get(0), L.get(1),L.get(2) );
        }
    }
);

JavaRDD<Row> employees = emps_f.map
(
    new Function< Tuple3<String, String, String> , Row>()
    {
        public Row  call(Tuple3<String, String, String> x)
        {
            return RowFactory.create( x._1(), Integer.parseInt(x._2()), Float.parseFloat( x._3()) );
        }
    }
);

StructType schema = new StructType(new StructField [] {
new StructField("name", DataTypes.StringType, false, Metadata.empty()),
new StructField("age", DataTypes.IntegerType, false, Metadata.empty()),
new StructField("salary", DataTypes.FloatType, false, Metadata.empty())
});

DataFrame people = sqlContext.createDataFrame(employees, schema);

people.show(3);

sqlContext.registerDataFrameAsTable(people,"employee");
DataFrame df2 = sqlContext.table("employee");

DataFrame sqltable = sqlContext.sql("select name from employee where salary >3500");

sqltable.show();

emps_f.saveAsTextFile(outputFile);

sc.stop();
}
}

```

Pom.xml

```

<project>
<groupId>hw6problem2</groupId>
<artifactId>spark-example</artifactId>
<modelVersion>4.0.0</modelVersion>
<name>problem2</name>
<packaging>jar</packaging>

```

```

<version>0.0.2</version>
<dependencies>
    <dependency> <!-- Spark dependency -->
        <groupId>org.apache.spark</groupId>
        <artifactId>spark-core_2.10</artifactId>
        <version>1.1.0</version>
        <scope>provided</scope>
    </dependency>
    <dependency>
        <groupId>org.apache.spark</groupId>
        <artifactId>spark-sql_2.10</artifactId>
        <version>1.5.1</version>
    </dependency>
</dependencies>
<properties>
    <java.version>1.8</java.version>
</properties>
<build>
    <pluginManagement>
        <plugins>
            <plugin>
                <groupId>org.apache.maven.plugins</groupId>
                <artifactId>maven-compiler-plugin</artifactId>
                <version>3.1</version>
                <configuration>
                    <source>${java.version}</source>
                    <target>${java.version}</target>
                </configuration>
            </plugin>
        </plugins>
    </pluginManagement>
</build>
</project>

```

Problem 3.

Attached file ebay.csv contains information of eBay's auction history. The Excel file has 9 columns and they represent:

The eBay online auction dataset has the following data fields:

- auctionid - unique identifier of an auction
- bid - the proxy bid placed by a bidder
- bidtime - the time (in days) that the bid was placed, from the start of the auction
- bidder - eBay username of the bidder
- bidderrate - eBay feedback rating of the bidder
- openbid - the opening bid set by the seller
- price - the closing price that the item sold for (equivalent to the second highest bid + an increment)
- item – name of the item being sold
- daystolive – length of the auction.

Using Spark DataFrames you will explore the data with following 4 questions:

1. How many auctions were held?
2. How many bids were made per item?
- 3.

- What's the minimum, maximum, and average bid (price) per item?
- What is the minimum, maximum and average number of bids per item?

4. Show the bids with price > 100

Import data into an RDD object. Transform that RDD into an RDD of Row-s by assign schema (column names and types). Transform that new RDD into a DataFrame. Call that DataFrame Auction.

```

CentOS_6.7_64-bit_2 - VMware Player (Non-commercial use only)
File Edit View Search Terminal Help
joe@localhost:~/Documents/Harvard63/HW6/hw6_problem3/src/main/java/ ... x
joe@localhost:~/Documents/Harvard63/HW6/hw6_problem3 ... x
File Edit View Search Terminal Help
joe@localhost:~/Documents/Harvard63/HW6/hw6_problem3/src/main/java/hw6problem3
[INFO] Nothing to compile - all classes are up to date
[INFO]
[INFO] --- maven-resources-plugin:2.6:testResources (default-testResources)
[WARNING] Using platform encoding (UTF-8 actually) to copy files.
[INFO] skip non existing resourceDirectory /home/joe/Documents/Harvard63/HW6/hw6_problem3/src/main/java/hw6problem3
[INFO] skip non existing resourceDirectory /home/joe/Documents/Harvard63/HW6/hw6_problem3/src/test/java
[INFO]
[INFO] --- maven-compiler-plugin:3.1:testCompile (default-testCompile)
[INFO] No sources to compile
[INFO]
[INFO] --- maven-surefire-plugin:2.12.4:test (default-test)
[INFO] No tests to run.
[INFO]
[INFO] --- maven-jar-plugin:2.4:jar (default-jar) @ spark-examples ...
[INFO] Building jar: /home/joe/Documents/Harvard63/HW6/hw6_problem3/target/spark-examples-1.0-SNAPSHOT.jar
[INFO] ...
[INFO] BUILD SUCCESS
[INFO] -----
[INFO] Total time: 12.466 s
[INFO] Finished at: 2016-03-11T05:13:01-08:00
[INFO] Final Memory: 14M/33M
[INFO] -----
joe@localhost hw6_problem3$ 

new StructField("bidder", DataTypes.StringType, false, Metadata.empty()),
new StructField("biderrate", DataTypes.IntegerType, false, Metadata.empty()),
new StructField("openbid", DataTypes.FloatType, false, Metadata.empty()),
new StructField("price", DataTypes.FloatType, false, Metadata.empty()),
new StructField("item", DataTypes.StringType, false, Metadata.empty()),
new StructField("daystolive", DataTypes.IntegerType, false, Metadata.empty())
];
DataFrame dfauction = sqlContext.createDataFrame(ebay_row, schema);

dfauction.show();
sqlContext.registerDataFrameAsTable(dfauction, "ebay");
DataFrame df_ebay2 = sqlContext.table("ebay");

DataFrame sqitable = sqlContext.sql("select * from ebay where bid >400");
sqitable.show();

//emps_f.saveAsTextFile(outputFile);
126,0-1

```

Intermediate output:

```

CentOS_6.7_64-bit_2 - VMware Player (Non-commercial use only)
File Edit View Search Terminal Help
Applications Places System joe@localhost:~/Documents/Harvard63/HW6/hw6_problem3/src/main/java/ - x
joe@localhost:~/Documents/Harvard63/HW6/hw6_problem3
File Edit View Search Terminal Help
joe@localhost:~/Documents/Harvard63/HW6/hw6_problem3/src/main/java/hw6problem3
File Edit View Search Terminal Help
joe@localhost:~/Documents/Harvard63/HW6/hw6_problem3/src/main/java/hw6problem3
[SLF4J: Found binding in [jar:file:/usr/lib/flume-ng/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class]
[SLF4J: Found binding in [jar:file:/usr/lib/zookeeper/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class]
[SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
[SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]
[auctionid| bid| bidtime| bidder|bidderrate
[8213034705| 95.0| 2.927373| jake7870|
[8213034705| 115.0| 2.943484| davidbresler2|
[8213034705| 100.0| 2.951285| gladimacowgirl|
only showing top 3 rows
[auctionid| bid| bidtime| bidder|bidderrate
[8212830525| 403.87| 2.152581| sybercool123|
[8212830525| 478.79| 2.588576| dqmignit|
[8212830525| 425.0| 4.471887| advisal|
[8212830525| 435.0| 4.858981| ruthhagerman|
[8212830525| 445.0| 4.859282| ruthhagerman|
[8212830525| 455.0| 4.859502| ruthhagerman|
[8212830525| 500.0| 4.859699| ruthhagerman|
"hw6p3.java" [dos] 154L, 5788C written

```

Running command

```
[joe@localhost hw6_problem2]$ spark-submit --class hw6problem3.hw6p3
./target/spark-example-0.0.2.jar
file:///home/joe/Documents/Harvard63/HW6/hw6_problem3/data/ebay.csv
```

Show (print) the schema of the DataFrame.

Relevant line of code

```
DataFrame dfauction = sqlContext.createDataFrame(ebay_row, schema);
dfauction.printSchema();
```

```
[joe@localhost hw6_problem3]$ spark-submit --class hw6problem3.hw6p3
./target/spark-example-0.0.2.jar
file:///home/joe/Documents/Harvard63/HW6/hw6_problem3/data/ebay.csv
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/lib/zookeeper/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/lib/flume-ng/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]
```

```
root
|-- auctionid: long (nullable = false)
|-- bid: float (nullable = false)
|-- bidtime: float (nullable = false)
|-- bidder: string (nullable = false)
|-- bidderrate: integer (nullable = false)
|-- openbid: float (nullable = false)
|-- price: float (nullable = false)
|-- item: string (nullable = false)
|-- daystolive: integer (nullable = false)
```

joe@localhost:~/Documents/Harvard63/HW6/hw6_problem3/src/main/java/ - □ x

```
[joe@localhost hw6_problem3]$ ls
[joe@localhost hw6_problem3]$ cd src/main/java/
[joe@localhost hw6_problem3]$ ls
[joe@localhost hw6_problem3]$ s
[joe@localhost hw6_problem3]$ spark-submit --class
file:///home/joe/Documents/Harvard63/HW6/hw6_pr
SLF4J: Class path contains multiple SLF4J bindin
SLF4J: Found binding in [jar:file:/usr/lib/zooke
icLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/lib/flume
cLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multi
SLF4J: Actual binding is of type [org.slf4j.impl
root
-- auctionid: long (nullable = false)
-- bid: float (nullable = false)
-- bidtime: float (nullable = false)
-- bidder: string (nullable = false)
-- bidderrate: integer (nullable = false)
-- openbid: float (nullable = false)
-- price: float (nullable = false)
-- item: string (nullable = false)
-- daystolive: integer (nullable = false)

[joe@localhost hw6_problem3]$
```

joe@localhost:~/Documents/Harvard63/HW6/hw6_problem3/src/main/java/hw6problem3 - □ x

```
[joe@localhost hw6_problem3]$ ls
[joe@localhost hw6_problem3]$ cd src/main/java/
[joe@localhost hw6problem3]$ ls
[joe@localhost hw6problem3]$ cat hw6p3.java
[joe@localhost hw6problem3]$
```

```
new StructField("price", DataTypes.FloatType, false, Metadata.empty()),
new StructField("item", DataTypes.StringType, false, Metadata.empty()),
new StructField("daystolive", DataTypes.IntegerType, false, Metadata.empty()))

};

DataFrame dfauction = sqlContext.createDataFrame(ebay_row, schema);

dfauction.printSchema();

//dfauction.show(3);

sqlContext.registerDataFrameAsTable(dfauction,"ebay");
DataFrame df_ebay2 = sqlContext.table("ebay");

DataFrame sqltable = sqlContext.sql("select * from ebay where bid >400");

//sqltable.show();

//emps_f.saveAsTextFile(outputFile);

sc.stop();
```

"hw6p3.java" [dos] 155L, 5824C written

149,7 99%

Make above queries using DataFrame API. You recall how we applied methods: select(), groupBy(), count() and others to the DataFrame in class. Use those methods.

1. How many auctions were held?

```
System.out.println( dfauction.select("auctionid").count());
```

```
joe@localhost:~/Documents/Harvard63/HW6/hw6_problem3/src/main/java/ - > joe@localhost:~/Documents/Harvard63/HW6/hw6_problem3/src/main/java/hw6problem3
```

```
[joe@localhost hw6_problem3]$ [joe@localhost hw6_problem3]$ spark-submit --cl [joe@localhost hw6_problem3]$ file:///home/joe/Documents/Harvard63/HW6/hw6_p [SLF4J: Class path contains multiple SLF4J bindin [SLF4J: Found binding in [jar:file:/usr/lib/zooke [SLF4J: Found binding in [jar:file:/usr/lib/flume [SLF4J: See http://www.slf4j.org/codes.html#multi [SLF4J: Actual binding is of type [org.slf4j.impl [root [-- auctionid: long (nullable = false) [-- bid: float (nullable = false) [-- bidtime: float (nullable = false) [-- bidder: string (nullable = false) [-- bidderrate: integer (nullable = false) [-- openbid: float (nullable = false) [-- price: float (nullable = false) [-- item: string (nullable = false) [-- daystolive: integer (nullable = false)
```

```
new StructField("openbid", DataTypes.FloatType, false, Metadata.empty()),  
new StructField("price", DataTypes.FloatType, false, Metadata.empty()),  
new StructField("item", DataTypes.StringType, false, Metadata.empty()),  
new StructField("daystolive", DataTypes.IntegerType, false, Metadata.empty())  
);  
  
DataFrame dfauction = sqlContext.createDataFrame(ebay_row, schema);  
  
dfauction.printSchema();  
  
System.out.println( dfauction.select("auctionid").count());  
  
//dfauction.show(3);  
//  
sqlContext.registerDataFrameAsTable(dfauction, "ebay");  
DataFrame df_ebay2 = sqlContext.table("ebay");  
  
DataFrame sqtable = sqlContext.sql("select * from ebay where bid >400");
```

```
10654 [joe@localhost hw6_problem3]$
```

```
"hw6p3.java" [dos] 160L, 5906C written
```

```
144,0-1 94
```

2. How many bids were made per item?

Relevant code line:

```
dfauction.groupBy("item").count().show();
```

Output:

```
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]  
root  
|-- auctionid: long (nullable = false)  
|-- bid: float (nullable = false)  
|-- bidtime: float (nullable = false)  
|-- bidder: string (nullable = false)  
|-- bidderrate: integer (nullable = false)  
|-- openbid: float (nullable = false)  
|-- price: float (nullable = false)  
|-- item: string (nullable = false)  
|-- daystolive: integer (nullable = false)  
  
+-----+----+  
| item|count|  
+-----+----+  
| xbox| 2784|  
| palm| 5917|  
| cartier| 1953|  
+-----+----+
```

```

[joe@localhost hw6_problem3]$ 
CentOS_6.7_64-bit_2 - VMware Player (Non-commercial use only)
Player | 
File Edit View Search Terminal Help
joe@localhost:~/Documents/Harvard63/HW6/hw6_problem3/src/main/java/ ... x
File Edit View Search Terminal Help
joe@localhost:~/Documents/Harvard63/HW6/hw6_problem3/src/main/java/ ... x
[ File Edit View Search Terminal Help
[ SLF4J: Found binding in [jar:file:/usr/lib/flume/cloggerBinder.class]
[ SLF4J: See http://www.slf4j.org/codes.html#multi
[ SLF4J: Actual binding is of type [org.slf4j.impl.root
[   |-- auctionid: long (nullable = false)
[   |-- bid: float (nullable = false)
[   |-- bidtime: float (nullable = false)
[   |-- bidder: string (nullable = false)
[   |-- bidderrate: integer (nullable = false)
[   |-- openbid: float (nullable = false)
[   |-- price: float (nullable = false)
[   |-- item: string (nullable = false)
[   |-- daystolive: integer (nullable = false)
[ +-----+-----+
[ |   item|count|
[ +-----+-----+
[ |   xbox| 2784|
[ |   palm| 5917|
[ |cartier| 1953|
[ +-----+-----+
[joe@localhost hw6_problem3]$ 
File Edit View Search Terminal Help
dfauction.printSchema();
//System.out.println("Total number of auctionid = "+ dfauction.select("auctionid").count());
dfauction.groupBy("item").count().show();
//dfauction.show(3);
//sqlContext.registerDataFrameAsTable(dfauction,"ebay");
DataFrame df_ebay2 = sqlContext.table("ebay");
DataFrame sqltable = sqlContext.sql("select * from ebay where bid > 0");
//sqltable.show();
//emps_f.saveAsTextFile(outputFile);
sc.stop();
"hw6p3.java" [dos] 164L, 5998C written

```

- What's the minimum, maximum, and average bid (price) per item?

```

dfauction.groupBy("item").min("price").show();
dfauction.groupBy("item").max("price").show();
dfauction.groupBy("item").mean("price").show();

```

```

+-----+-----+
| item|min(price)| 
+-----+-----+
|  xbox| 31.0|
|  palm| 175.0|
|cartier| 26.0|
+-----+-----+ 

+-----+-----+
| item|max(price)| 
+-----+-----+
|  xbox| 501.77|
|  palm| 290.0|
|cartier| 5400.0|
+-----+-----+ 

+-----+-----+
| item| avg(price)| 
+-----+-----+
|  xbox| 144.2759409416681|
|  palm| 231.13026672693368|
|cartier| 925.0479063882744|
+-----+-----+

```

```

+-----+
| item[min(price)]|
+-----+
| xbox| 31.0|
| palm| 175.0|
| cartier| 26.0|
+-----+-----+
+-----+
| item[max(price)]|
+-----+
| xbox| 501.77|
| palm| 290.0|
| cartier| 5400.0|
+-----+-----+
+-----+
| item| avg(price)|
+-----+
| xbox| 144.2759409416681|
| palm| 231.13026672693368|
| cartier| 925.0479663882744|
+-----+-----+
[joe@localhost hw6_problem3]$ ./hw6p3.java [dos] 169L, 6162C written

```

What is the minimum, maximum and average **number of bids** per item?

Relevant code: the first line counts the number of bids per item. The second line applies operation describe() to that result. The first line is given only to understand better the outcome of the second line, which is the only relevant to the posed question.

```
dfauction.groupBy("item").count().show();
```

```
dfauction.groupBy("item").count().describe().show();
```

Output:

```
dfauction.groupBy("item").count().show();
```

```
+-----+
| item|count| the "count" counts the number of bids per item
+-----+
| xbox| 2784|
| palm| 5917|
| cartier| 1953|
-----+
```

```
dfauction.groupBy("item").count().describe().show();
```

```
+-----+
|summary| count|
+-----+
| count| 3 |
| mean| 3551.333333333335 |
| stddev| 1706.834171467424 |
| min| 1953 |
| max| 5917 |
-----+
```

The screenshot shows a dual-terminal session on a CentOS 6.7 system. The left terminal window displays the contents of `hw6problem3.java`, which includes code for calculating the total number of auction IDs, minimum, maximum, and average bid per item, and the count of bids per item. The right terminal window shows the execution of the Java code, resulting in the output:

```
joe@localhost:~/Documents/Harvard63/HW6/hw6_problem3/src/main/java/hw6problem3
```

```
//System.out.println("Total number of auctionid = "+dfauction.select("auctionid").count());  
dfauction.groupBy("item").count().show();  
// What's the minimum, maximum, and average bid (price) per item?  
dfauction.groupBy("item").min("price").show();  
dfauction.groupBy("item").max("price").show();  
dfauction.groupBy("item").mean("price").show();  
  
// What is the minimum, maximum and average number of bids per item?  
//dfauction.groupBy().min().describe("item").show();  
dfauction.groupBy("item").count().describe().show();  
  
//dfauction.groupBy("item").max("price").show();  
//dfauction.groupBy("item").mean("price").show();  
  
"hw6p3.java" [dos] 181L, 6514C written
```

The status bar at the bottom indicates the session is at line 157, offset 0-1, with 89% completion.

Show the bids with price > 100

```
dfauction.filter("price>100").select("bid").show(10);
```

Extension

```
dfauction.filter("price>100").select("bid","price","item").show(10);
```

```
joe@localhost:~/Documents/Harvard63/HW6/hw6_problem3/src/main/java/i_ - x
joe@localhost:~/Documents/Harvard63/HW6/hw6_problem3 - x
joe@localhost:~/Documents/Harvard63/HW6/hw6_problem3$ - x
File Edit View Search Terminal Help
File Edit View Search Terminal Help
File Edit View Search Terminal Help
//System.out.println("Total number of auctionid = "+ dfauction.select("auctionid").count());
dfauction.groupBy("item").count().show();
// What's the minimum, maximum, and average bid (price) per item?
// dfauction.groupBy("item").min("price").show();
// dfauction.groupBy("item").max("price").show();
// dfauction.groupBy("item").mean("price").show();
// What is the minimum, maximum and average number of bids per item?
// dfauction.groupBy("item").count().describe().show();
//Show the bids with price > 100
dfauction.filter("price>100").select("bid","price","item").show(10);
only showing top 10 rows
"hw6p3.java" [dos] 185L, 6577C written
161,0-1 87%
[ctrl] [joe@localhost hw6_problem3]$
```

Next transform your Auction DataFrame into a table

Transform into a table and test it with a select query. Relevant to this question lines are bolded.

```
sqlContext.registerDataFrameAsTable(dfauction,"ebay");
DataFrame df_ebay = sqlContext.table("ebay");
DataFrame sqltable = sqlContext.sql("select * from ebay where price >200 ");
sqltable.show(10);
```

The output:

CentOS_6.7_64-bit_2 - VMware Player (Non-commercial use only)

Player Applications Places System joe@localhost Fri Mar 11, 11:40 AM

joe@localhost:~/Documents/Harvard63/HW6/hw6_problem3

```
[File Edit View Search Terminal Help]
[-- daystolive: integer (nullable = false)
[| item:|--+
[| | item|count|
[| |+-----+
[| | | xbox| 2784|
[| | | palm| 5917|
[| | cartier| 1953|
[| +-----+
[| h+-----+
[| | auctionid| bid| bidtime| bidder| bidderate| openbid|
[| +-----+
[| | 8213935134| 20.0| 0.799213| arlysanhawk| 9| 0.1|
[| /8213935134| 1.0| 0.83419| isa505| 1| 0.1|
[| | 8213935134| 5.0| 0.837257| isa505| 1| 0.1|
[| | 8213935134| 10.0| 1.018333| bluebubbles| 1| 26| 0.1|
[| | 8213935134| 15.0| 1.01838| bluebubbles| 1| 26| 0.1|
[| | 8213935134| 22.01| 1.018438| bluebubbles| 1| 26| 0.1|
[| | 8213935134| 55.0| 1.153345| drew6330| 19| 0.1|
[| | 8213935134| 50.0| 1.307049| keinpasty| 11| 0.1|
[| | 8213935134| 60.0| 1.307222| keinpasty| 11| 0.1|
[| | 8213935134| 105.51| 1.639387| txmp| 28| 0.1|
[| +-----+
only showing top 10 rows
```

[joe@localhost hw6_problem3]\$

joe@localhost:~/Documents/Harvard63/HW6/hw6_problem3/src/main/java/hw6problem3

```
[File Edit View Search Terminal Help]
// dfaction.groupby("item").count().describe().show();

//Show the bids with price > 100
//OK    dfaction.filter("price>100").select("bid","price","item").show(10);

//Next transform your Auction DataFrame into a table

sqlContext.registerDataFrameAsTable(dfaction,"ebay");
DataFrame df_ebay = sqlContext.table("ebay");

DataFrame sqltable = sqlContext.sql("select * from ebay where price >200 ");

sqltable.show(10);

//emps_f.saveAsTextFile(outputFile);

sc.stop();
```

"hw6p3.java" [dos] 177L, 6511C written

159,

Make the same 4 inquiries using regular SQL queries.

```
// how many auctions  
sqlContext.sql("select count(auctionid) from ebay ").show(10);
```

```
// how many bids  
sqlContext.sql("select item, count(*) from ebay group by item").show(10);
```

CentOS_6.7_64-bit_2 - VMware Player (Non-commercial use only)

Player Applications Places System joe Fri Mar 11, 1:08 PM

```
joe@localhost:~/Documents/Harvard63/HW6/hw6_problem3
joe@localhost:~/Documents/Harvard63/HW6/hw6_problem3/src/main/java/hw6problem3
```

f File Edit View Search Terminal

```
[| 8213935134| 50.0|1.307049]
[| 8213935134| 60.0|1.307222]
[| 8213935134| 105.51|1.639387]
```

//only showing top 10 rows

//Next transform your Auction DataFrame into a table

```
sqlContext.registerDataFrameAsTable(dfauction,"ebay");
DataFrame df_ebay = sqlContext.table("ebay");

DataFrame sqtable = sqlContext.sql("select * from ebay where price >200 ");
sqtable.show(10);

// how many auctions
sqlContext.sql("select count(auctionid) from ebay ").show(10);

// how many bids
sqlContext.sql("select item, count(*) from ebay group by item").show(10);

// What's the minimum, maximum, and average bid (price) per item?
sqlContext.sql("select item, min(price) from ebay group by item").show(10);
sqlContext.sql("select item, max(price) from ebay group by item").show(10);
sqlContext.sql("select item, avg(price) from ebay group by item").show(10);

-- INSERT --
```

158,1 87%

joe@localhost:/mnt/h... joe@localhost:~/Docu... joe@localhost:~/Docu... joe@localhost:~/Docu...

```
// What's the minimum, maximum, and average bid (price) per item?  
sqlContext.sql("select min(price) from ebay group by item").show(10);  
sqlContext.sql("select max(price) from ebay group by item").show(10);  
sqlContext.sql("select mean(price) from ebay group by item").show(10);
```

```
//What is the minimum, maximum and average number of bids per item?
sqlContext.sql("select min(cnt) from (select item, count(*) as cnt from ebay group by item) as t1").show(10);
sqlContext.sql("select max(cnt) from (select item, count(*) as cnt from ebay group by item) as t1").show(10);
sqlContext.sql("select avg(cnt) from (select item, count(*) as cnt from ebay group by item) as t1").show(10);
```

```
joe@localhost:~/Documents/Harvard63/HW6/hw6_problem3
File Edit View Terminal Help
DataFrame sqltable = sqlContext.sql("select * from ebay where price >200 ");
sqltable.show(10);

// how many auctions
sqlContext.sql("select count(auctionid) from ebay ").show(10);

// how many bids
sqlContext.sql("select item, count(*) from ebay group by item").show(10);

// What's the minimum, maximum, and average bid (price) per item?
sqlContext.sql("select item, min(price) from ebay group by item").show(10);
sqlContext.sql("select item, max(price) from ebay group by item").show(10);
sqlContext.sql("select item, avg(price) from ebay group by item").show(10);

//What is the minimum, maximum and average number of bids per item?
sqlContext.sql("select min(cnt) from (select item, count(*) as cnt from ebay group by item) as t1").show(10);
sqlContext.sql("select max(cnt) from (select item, count(*) as cnt from ebay group by item) as t1").show(10);
sqlContext.sql("select avg(cnt) from (select item, count(*) as cnt from ebay group by item) as t1").show(10);

-- INSERT --
```

```
//Show the bids with price > 100
sqlContext.sql("select bid from ebay where price>100").show(10);
```

```
joe@localhost:~/Documents/Harvard63/HW6/hw6_problem3
File Edit View Search Terminal Help
sqlContext.sql("select min(cnt) from (select item, count(*) as cnt from ebay group by item) as t1").show(10);
sqlContext.sql("select max(cnt) from (select item, count(*) as cnt from ebay group by item) as t1").show(10);
sqlContext.sql("select avg(cnt) from (select item, count(*) as cnt from ebay group by item) as t1").show(10);

//Show the bids with price > 100
sqlContext.sql("select bid from ebay where price>100").show(10);

//chemaPeople.write().parquet("people.parquet");

only showing//sqltable.show(10);
//
//emps.f.saveAsTextFile(outputFile);

sc.stop();
```

Persist your DataFrame as a Parquet file and show that you could exit your pyspark shell and come back in it and you will be still able to read the data from that file and create a DataFrame and an SQL like table that you can issue queries against.

```
//parquet
    sqltable.write().parquet("dfebay.parquet");

// Read in the parquet file created above.
DataFrame parquetFile = sqlContext.read().parquet("dfebay.parquet");

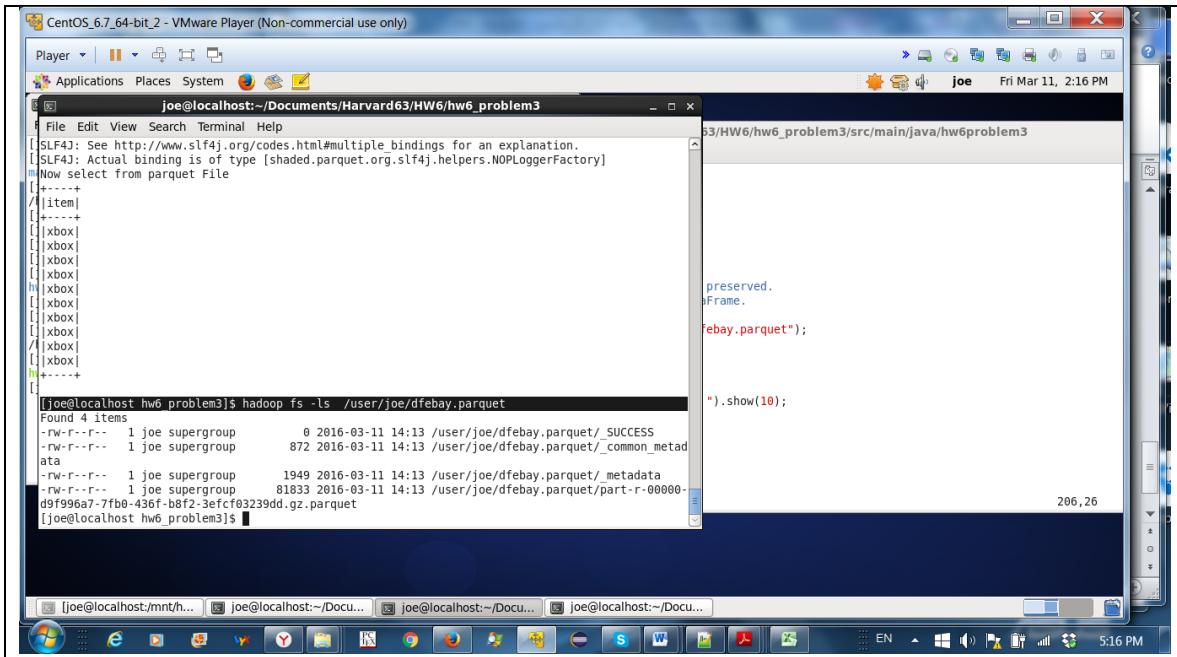
parquetFile.registerTempTable("parquetFile");
System.out.println("Now select from parquet File");
sqlContext.sql("SELECT item FROM parquetFile LIMIT 10").show(10);
```

The screenshot shows a terminal window titled "CentOS_6.7-64-bit_2 - VMware Player (Non-commercial use only)". The window contains two tabs: "joe@localhost:~/Documents/Harvard63/HW6/hw6_problem3" and "joe@localhost:~/Documents/Harvard63/HW6/hw6_problem3/src/main/java/hw6problem3". The code in the tabs is identical to the one provided in the previous code block. The output shows the logs of the Java application, which includes SLF4J binding information and the execution of the Scala code to write and read the Parquet file. The terminal also shows the resulting DataFrame being registered as a temporary table and the execution of an SQL query to select the first 10 items.

The parquet file is located in the hadoop hdfs so it certainly exists beyond the life-span of the screen I work in.

`hadoop fs -ls /user/joe/dfebay.parquet`

Found 4 items		
-rw-r--r--	1 joe supergroup	0 2016-03-11 14:13 /user/joe/dfebay.parquet/_SUCCESS
-rw-r--r--	1 joe supergroup	872 2016-03-11 14:13 /user/joe/dfebay.parquet/_common_metadata
-rw-r--r--	1 joe supergroup	1949 2016-03-11 14:13 /user/joe/dfebay.parquet/_metadata
-rw-r--r--	1 joe supergroup	81833 2016-03-11 14:13 /user/joe/dfebay.parquet/part-r-00000-d9f996a7-7fb0-436fb8f2-3efcf03239dd.gz.parquet



Entire output:

```
joe@localhost hw6_problem3]$ spark-submit --class hw6problem3.hw6p3 ./target/spark-example-0.0.2.jar
file:///home/joe/Documents/Harvard63/HW6/hw6_problem3/data/ebay.csv
^C[joe@localhost hw6_problem3]$  

[joe@localhost hw6_problem3]$ hadoop fs -rm -r /user/joe/dfebay.parquet
16/03/11 14:40:28 INFO fs.TrashPolicyDefault: Namenode trash configuration: Deletion interval = 0
minutes, Emptier interval = 0 minutes.
Deleted /user/joe/dfebay.parquet
[joe@localhost hw6_problem3]$ spark-submit --class hw6problem3.hw6p3 ./target/spark-example-0.0.2.jar
file:///home/joe/Documents/Harvard63/HW6/hw6_problem3/data/ebay.csv
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/lib/zookeeper/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/lib/flume-ng/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]
root
|- auctionid: long (nullable = false)
|- bid: float (nullable = false)
|- bidtime: float (nullable = false)
|- bidder: string (nullable = false)
|- biderrate: integer (nullable = false)
|- openbid: float (nullable = false)
|- price: float (nullable = false)
|- item: string (nullable = false)
|- daystolive: integer (nullable = false)

Total number of auctionid = 10654
+-----+
| item|count|
+-----+
| xbox| 2784|
| palm| 5917|
| cartier| 1953|
+-----+

+-----+
| item|min(price)|
+-----+
| xbox| 31.0|
| palm| 175.0|
| cartier| 26.0|
+-----+

+-----+
| item|max(price)|
+-----+
| xbox| 501.77|
```

```

| palm| 290.0|
| cartier| 5400.0|
+-----+-----+
+-----+-----+
| item| avg(price)|
+-----+-----+
| xbox| 144.2759409416681|
| palm|231.13026672693368|
| cartier| 925.0479063882744|
+-----+-----+
+-----+-----+
|summary| count|
+-----+-----+
| count| 3|
| mean|3551.333333333335|
| stddev| 1706.834171467424|
| min| 1953|
| max| 5917|
+-----+-----+
+-----+-----+
| bid|price|item|
+-----+-----+----+
| 95.0|117.5|xbox|
|115.0|117.5|xbox|
|100.0|117.5|xbox|
|117.5|117.5|xbox|
| 2.0|120.0|xbox|
|15.25|120.0|xbox|
| 3.0|120.0|xbox|
|10.0|120.0|xbox|
|24.99|120.0|xbox|
| 20.0|120.0|xbox|
+-----+-----+
only showing top 10 rows

+-----+-----+-----+-----+-----+-----+-----+
| auctionid| bid| bidtime| bidder|bidderrate|openbid|price|item|daystolive|
+-----+-----+-----+-----+-----+-----+-----+
|8213935134| 20.0|0.799213| arlysanhawk| 9| 0.1|207.5|xbox| 3|
|8213935134| 1.0| 0.83419| isai505| 1| 0.1|207.5|xbox| 3|
|8213935134| 5.0|0.837257| isai505| 1| 0.1|207.5|xbox| 3|
|8213935134| 10.0|1.018333|bluebubbles_1| 26| 0.1|207.5|xbox| 3|
|8213935134| 15.0| 1.01838|bluebubbles_1| 26| 0.1|207.5|xbox| 3|
|8213935134| 22.01|1.018438|bluebubbles_1| 26| 0.1|207.5|xbox| 3|
|8213935134| 55.0|1.153345| drew6330| 19| 0.1|207.5|xbox| 3|
|8213935134| 50.0|1.307049| keipnasty| 11| 0.1|207.5|xbox| 3|
|8213935134| 60.0|1.307222| keipnasty| 11| 0.1|207.5|xbox| 3|
|8213935134|105.51|1.639387| txmp| 28| 0.1|207.5|xbox| 3|
+-----+-----+-----+-----+-----+-----+
only showing top 10 rows

+----+
| _c0|
+----+
|10654|
+----+
+----+----+
| item| _c1|
+----+----+
| xbox|2784|
| palm|5917|
| cartier|1953|
+----+----+
+----+----+
| item| _c1|
+----+----+
| xbox| 31.0|
| palm|175.0|
| cartier| 26.0|
+----+----+
+----+----+
| item| _c1|
+----+----+
| xbox|501.77|
| palm| 290.0|

```

```
|cartier|5400.0|
+-----+
+-----+-----+
| item|      _c1|
+-----+-----+
| xbox| 144.2759409416681|
| palm|231.13026672693368|
|cartier| 925.0479063882744|
+-----+-----+


+---+
| _c0|
+---+
|1953|
+---+


+---+
| _c0|
+---+
|5917|
+---+

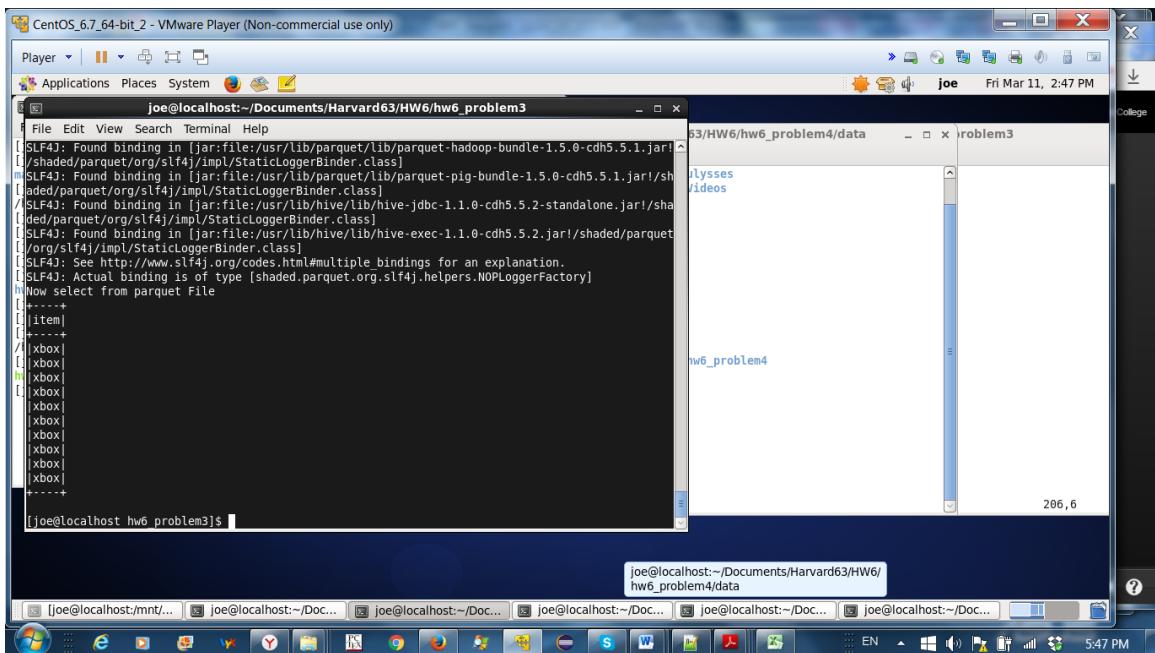

+-----+
|       _c0|
+-----+
|3551.33333333333335|
+-----+


+---+
| bid|
+---+
| 95.0|
|115.0|
|100.0|
|117.5|
| 2.0|
|15.25|
| 3.0|
| 10.0|
|24.99|
| 20.0|
+---+
only showing top 10 rows

SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/lib/parquet/lib/parquet-format-2.1.0-cdh5.5.1.jar!/:shaded/parquet/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/lib/parquet/lib/parquet-hadoop-bundle-1.5.0-cdh5.5.1.jar!/:shaded/parquet/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/lib/parquet/lib/parquet-pig-bundle-1.5.0-cdh5.5.1.jar!/:shaded/parquet/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/lib/hive/lib/hive-jdbc-1.1.0-cdh5.5.2-standalone.jar!/:shaded/parquet/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/lib/hive/lib/hive-exec-1.1.0-cdh5.5.2.jar!/:shaded/parquet/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [shaded.parquet.org.slf4j.helpers.NOPLoggerFactory]

Now select from parquet File
+---+
|item|
+---+
|xbox|
|xbox|
|xbox|
|xbox|
|xbox|
|xbox|
|xbox|
|xbox|
|xbox|
+---+


[joe@localhost hw6_problem3]$
```



The code:

```
/*
 * @author Serguey Khovansky
 * @version: March 2016
 * This is solution of the homework 6, problem 3:
 *
 * Attached file ebay.csv contains information of eBay's auction history.
 * The Excel file has 9 columns and they represent:
 * The eBay online auction dataset has the following data fields:
 * • auctionid - unique identifier of an auction
 * • bid - the proxy bid placed by a bidder
 * • bidtime - the time (in days) that the bid was placed, from the start of the auction
 * • bidder - eBay username of the bidder
 * • bidderrate - eBay feedback rating of the bidder
 * • openbid - the opening bid set by the seller
 * • price - the closing price that the item sold for (equivalent to
 *           the second highest bid + an increment)
 * • item - name of the item being sold
 * • daystolive - length of the auction.
 *
 * Using Spark DataFrames you will explore the data with following 4 questions:
 * 1. How many auctions were held?
 * 2. How many bids were made per item?
 * 3.
 * o What's the minimum, maximum, and average bid (price) per item?
 * o What is the minimum, maximum and average number of bids per item?
 * 4. Show the bids with price > 100
 * Import data into an RDD object. Transform that RDD into an RDD of Row-s by assign schema
```

```

* (column names and types). Transform that new RDD into a DataFrame. Call that DataFrame Auction.
* Show (print) the schema of the DataFrame. Make above queries using DataFrame API.
* You recall how we applied methods: select(), groupBy(), count() and others to the DataFrame
* in class. Use those methods.
* Next transform your Auction DataFrame into a table and make the same 4 inquiries
* using regular SQL queries.
*
*/

```

```

package hw6problem3;

import java.util.Arrays;
import java.util.List;
import java.lang.Iterable;

import scala.Tuple9;

import org.apache.commons.lang.StringUtils;

import org.apache.spark.SparkConf;
import org.apache.spark.api.java.JavaRDD;
import org.apache.spark.api.java.JavaPairRDD;
import org.apache.spark.api.java.JavaSparkContext;
import org.apache.spark.api.java.function.FlatMapFunction;
import org.apache.spark.api.java.function.Function;

import org.apache.spark.sql.SQLContext;
import org.apache.spark.sql.Row;
import org.apache.spark.sql.RowFactory;

import org.apache.spark.sql.types.*;
import org.apache.spark.sql.DataFrame;

/**
 * The class hw6p1 make database operations with two text files using SPARK MapReduce technology
 */
public class hw6p3 {

    /**
     * @param args[] input and output paths and location of the file input and output files
     */
    public static void main(String[] args) throws Exception {
        String inputFile = args[0];
        // String outputFile= args[1];

        /**
         * Create a Java Spark Context sc
         */

```

```

SparkConf conf = new SparkConf().setAppName("hw6p3");
JavaSparkContext sc = new JavaSparkContext(conf);
SQLContext sqlContext = new SQLContext(sc);

/***
 *      Create RDD emps
 */
JavaRDD<String> ebay = sc.textFile(inputFile);

/***
 *Turn into Tuple of words
*/
JavaRDD<Tuple9<String, String, String, String, String, String, String, String, String>> ebay_rdd_tuple9 =
ebay.map
(
    new Function<String, Tuple9<String, String, String, String, String, String, String, String, String>>()
    {
        public Tuple9<String, String, String, String, String, String, String, String, String> call(String x)
        {
            List L= Arrays.asList(x.toLowerCase().split( "\\", ) );
            return new Tuple9<String, String, String, String, String, String, String, String, String>(L.get(0), L.get(1),L.get(2), L.get(3),L.get(4),L.get(5),L.get(6),L.get(7),L.get(8) );
        }
    }
);

JavaRDD<Row> ebay_row = ebay_rdd_tuple9.map
(
    new Function< Tuple9<String, String, String, String, String, String, String, String, String> , Row>()
    {
        public Row call(Tuple9<String, String, String, String, String, String, String, String, String> x)
        {
            return RowFactory.create( Long.parseLong(x._1()), Float.parseFloat( x._2() ),Float.parseFloat(x._3()), x._4(),
Integer.parseInt(x._5()), Float.parseFloat( x._6()),Float.parseFloat( x._7()), x._8().toString(),Integer.parseInt( x._9() ) );
        }
    }
);

StructType schema = new StructType(new StructField [] {
new StructField("auctionid", DataTypes.LongType, false, Metadata.empty()),
new StructField("bid", DataTypes.FloatType, false, Metadata.empty()),

```

```

new StructField("bidtime", DataTypes.FloatType, false, Metadata.empty()),
new StructField("bidder", DataTypes.StringType, false, Metadata.empty()),
new StructField("bidderrate", DataTypes.IntegerType, false, Metadata.empty()),
new StructField("openbid", DataTypes.FloatType, false, Metadata.empty()),
new StructField("price", DataTypes.FloatType, false, Metadata.empty()),
new StructField("item", DataTypes.StringType, false, Metadata.empty()),
new StructField("daystolive", DataTypes.IntegerType, false, Metadata.empty())

});

DataFrame dfauction = sqlContext.createDataFrame(ebay_row, schema);

dfauction.printSchema();

System.out.println("Total number of auctionid = "+ dfauction.select("auctionid").count());

dfauction.groupBy("item").count().show();

// What's the minimum, maximum, and average bid (price) per item?
dfauction.groupBy("item").min("price").show();
dfauction.groupBy("item").max("price").show();
dfauction.groupBy("item").mean("price").show();

// What is the minimum, maximum and average number of bids per item?

dfauction.groupBy("item").count().describe().show();

//Show the bids with price > 100

dfauction.filter("price>100").select("bid","price","item").show(10);

//Next transform your Auction DataFrame into a table

sqlContext.registerDataFrameAsTable(dfauction,"ebay");
DataFrame df_ebay = sqlContext.table("ebay");

DataFrame sqltable = sqlContext.sql("select * from ebay where price >200 ");
sqltable.show(10);

// how many auctions
sqlContext.sql("select count(auctionid) from ebay ").show(10);

// how many bids
sqlContext.sql("select item, count(*) from ebay group by item").show(10);

```

```

// What's the minimum, maximum, and average bid (price) per item?
sqlContext.sql("select item, min(price) from ebay group by item").show(10);
sqlContext.sql("select item, max(price) from ebay group by item").show(10);
sqlContext.sql("select item, avg(price) from ebay group by item").show(10);

//What is the minimum, maximum and average number of bids per item?
sqlContext.sql("select min(cnt) from (select item, count(*) as cnt from ebay group by item) as t1").show(10);
sqlContext.sql("select max(cnt) from (select item, count(*) as cnt from ebay group by item) as t1").show(10);
sqlContext.sql("select avg(cnt) from (select item, count(*) as cnt from ebay group by item) as t1").show(10);

//Show the bids with price > 100
sqlContext.sql("select bid from ebay where price>100").show(10);

//parquet
sqltable.write().parquet("dfebay.parquet");

// Read in the parquet file created above.
// Parquet files are self-describing so the schema is preserved.
// The result of loading a parquet file is also a DataFrame.

DataFrame parquetFile = sqlContext.read().parquet("dfebay.parquet");

parquetFile.registerTempTable("parquetFile");
// Dataset<Row> tmp =
System.out.println("Now select from parquet File");
sqlContext.sql("SELECT item FROM parquetFile LIMIT 10 ").show(10);

sc.stop();
}
}

```

POM.xml

```

<project>
<groupId>hw6problem3</groupId>
<artifactId>spark-example</artifactId>
<modelVersion>4.0.0</modelVersion>
<name>problem3</name>
<packaging>jar</packaging>
<version>0.0.2</version>
<dependencies>
<dependency> <!-- Spark dependency -->
<groupId>org.apache.spark</groupId>
<artifactId>spark-core_2.10</artifactId>
<version>1.1.0</version>
<scope>provided</scope>
</dependency>
<dependency>

```

```

<groupId>org.apache.spark</groupId>
<artifactId>spark-sql_2.10</artifactId>
<version>1.5.1</version>
</dependency>
<dependency>
<groupId>com.twitter</groupId>
<artifactId>parquet-avro</artifactId>
<version>1.0.0</version>
</dependency>
</dependencies>
<properties>
<java.version>1.8</java.version>
</properties>
<build>
<pluginManagement>
<plugins>
<plugin>
<groupId>org.apache.maven.plugins</groupId>
<artifactId>maven-compiler-plugin</artifactId>
<version>3.1</version>
<configuration>
<source>${java.version}</source>
<target>${java.version}</target>
</configuration>
</plugin>
</plugins>
</pluginManagement>
</build>
</project>

```

Problem 4.

Use Sqoop to import all tables in MySQL demo database `retail_db` into Hive.
Use Spark DataFrame API or Spark Temporary Tables to find orders with the largest number of order items per order.

To make `retail_db` import into Hive

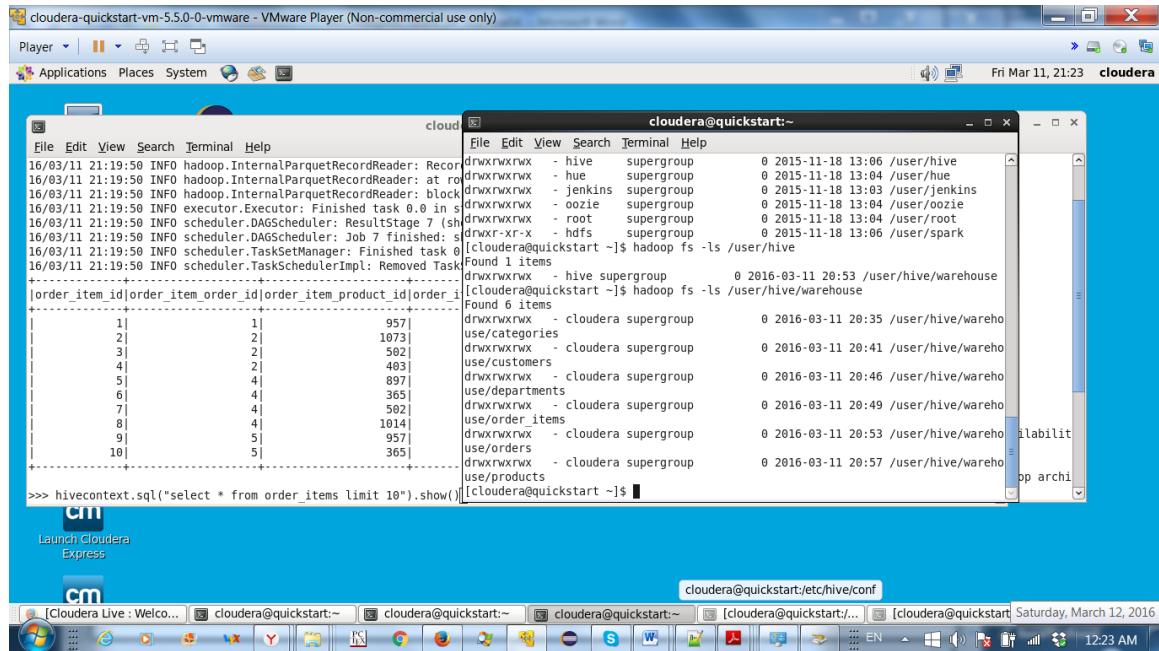
```

sqoop import-all-tables -m 1 --connect jdbc:mysql://quickstart:3306/retail_db --
username=retail_dba --password=cloudera --compression-codec=snappy --as-parquetfile
--warehouse-dir=/user/hive/warehouse --hive-import --hive-overwrite

```

P.S. Could not make Sqoop work as required for more than 7 hours....It began working after re-installing VM.

Required files of the database are in the Hadoop along with the results of the test query



```
// Note, that sc is produced by pyspark by default
from pyspark import SparkConf, SparkContext
hivecontext = HiveContext(sc)
```

```
//test query
hivecontext.sql("select * from order_items limit 10").show()
```

Find orders with the largest number of order items per order

```
hivecontext.sql("select max(oi) from (select order_item_quantity as oi from order_items ) as t1").show(10);
```

Output:

```
+---+
| _c0 |
+---+
| 5 |
+---+
```

