**Project 7, Big Data Analysis**

Kristie Kooken

Due: March 2, 2024

For project 7, Google Cloud was utilized to create a VM instance and run analyses. Project 7 conducted all analyses in this environment. I investigated data from E-Commerce - Product Inventory Management (www.kaggle.com) to gain insights into sales trends and reordering points. The e-commerce dataset has about 500,000 records with 8 column variables: Invoice Number as a numeric variable, Stock Code as a categorical variable, Description as a text character description, Quantity as a numeric variable (integer) for the number of items in the Description, Invoice Date as a character variable (date of transaction), Unit Price as a numeric variable representing the cost of the item in the Description, Customer ID as a numeric variable representing the customer's identifier and Country as a text variable representing the country selling the item in the Description.

Specific questions of interest included:

- What are the top 5 countries that have the highest prices?
- What are the top 5 countries that have the lowest prices?
- Who are the top 5 performing countries (sells the most quantities) in the world?
- Who are the lowest performing countries (sells the least quantities) in the world?
- Does pricing predict quantity?

To investigate these questions, this data was ingested in NiFi and analyzed using Spark, Pyspark. Initially, I did use scala to read the .json format to a dataframe and I describe challenges I encountered with this process below. I wanted to bring the data in using NiFi as I was interested to learn more about how NiFi might work to ingest larger amounts of

data in what I thought might be a more common data flow. Though I did a smaller amount of analysis in scala – I ended up redoing all analysis in Pyspark. I picked Pyspark as it has many powerful components and is more intuitive from my perspective. As it is very widely used, I wanted to gain more experience with it.

To bring the data into NiFi, I copied the e-commerce dataset to our Google Cloud VM location using scp and copied the data to the appropriate folder location for NiFi. Then, I created a flow in NiFi to read the .csv into NiFi and create a .json file. Configuring the NiFi flow was challenging and after following a few different methods, I decided to use a tutorial and template as a guide (www.community.cloudera.com). Screenshot 1 shows the successful NiFi flow for this process and screenshot 2 shows the original file and the converted .json format of the file saved to the /dsc650-infra/bellevue-bigdata/hadoop-hive-spark-hbase/data location. Screenshots 3 and 4 show contents of each file in the original and json format.

Once the data was uploaded to the data folder, I read the file into a dataframe using Spark scala (Screenshot 5). I encountered many challenges doing this, and then learned that json can be 1.5 to 3 times bigger than .csv format (www.jsoneditoronline.org). A file of this size was very challenging to process on our VM instance (Screenshot 6) where our cluster does not have a lot of compute and because of this, I decided to move forward with analyzing the .csv file instead.

Because I moved forward with analyzing the .csv format, I decided to use Pyspark for data cleaning and analyses. After initially reading in the data, I looked at the highest or lowest values of frequency counts as well as the country list to determine if there were values that should be removed from analysis. Data cleaning steps included filtering the

Country variable to remove "Unspecified" as well as "European Community" and only

keep those values of Unit Price that were greater than 0.10. Also, any missing values as

well as any duplicates were removed from the data. Screenshot 7 shows reading in the

data and creating a dataframe and changing numeric strings to numeric columns.

Screenshot 8a shows frequencies of each column and displays the head or the tail to

inspect each column for values that are not allowed for analysis (e.g., a negative value for

Unit Price). Screenshot 8b shows the frequency of the column Country to clean this

variable of entries that do not represent a single country. Screenshot 8c shows data

cleaning steps of filtering and then removing duplicate rows as well as missing values.

Screenshot 9 is the first 20 rows of the dataset used for analysis.

To investigate the highest and lowest performing countries in terms of item prices

and quantities of items, a grouped frequency that sorted the data from highest to lowest

was run, Screenshot 10 has the results for quantity and Screenshot 11 has the results for

pricing.

Initially I was interested in running correlations or descriptive statistics on the

dataframe and used the Summarizer package in Pyspark to investigate running these

analyses. However, I was not able to run this error-free and need to further investigate

what needs to be corrected with my code. Screenshot 12 shows the test data run I

conducted in order to better understand the input format and vectors

(www.spark.apache.org).

The final analysis that was run in Pyspark was a linear regression using Quantity

as the dependent variable and Price as the independent variable, Screenshot 13. This

simple regression model did not yield good results, Screenshot 14. I was very interested

to get to this stage and explore how to run a model in Pyspark. There is a lot more to learn but it is interesting and motivating to scratch the surface of these analyses.

There were different issues that I ran into while completing this project, some I figured out and some are still being worked on. An issue that caused a lot of re-work was not knowing how to install packages of python that were missing (e.g., numpy) and though I figured this out (install on all nodes), it made me realize that it is quite complex to have a production environment even with installing a new package. Likewise, I spent a long time trying to install a plotting package until I read several articles on how plotting is often done outside of a big data architecture.

In closing, this project allowed me to explore a variety of tools and get a better understanding of how big data architecture works from data flow to analyses. There is a lot more to learn, however I enjoyed learning more about NiFi, scala and Pyspark as mentioned here.

Screenshots referenced in this document:

#1:

#2:



#3:

#4:



#5:

#6:



#7:

#8a:



#8b:

#8c:



```
>>> df2 = spark.sql("SELECT * FROM df WHERE UnitPriceN > .10 and Country != 'Unspecified' and Country != 'European Community' and CustomerID is not null")
>>> df.createOrReplaceTempView("df2")
>>> from pyspark.sql.functions import col
>>> print("Distinct count: "+str(df2.count()))
Distinct count: 406177
>>> df2_d = df2.distinct()
>>> print("Distinct count: "+str(df2_d.count()))
Distinct count: 400958
>>> df.createOrReplaceTempView("df2_d")
>>> df_any = df2_d.dropna(how="any")
>>> print("Distinct count: "+str(df_any.count()))
Distinct count: 392095
>>> df.createOrReplaceTempView("df_any")
>>>
```

#9:



```
>>> print("Distinct count: "+str(df_any.count()))
Distinct count: 392095
>>> df_any.show()
+---------+---------+--------------------+--------+--------------+---------+--------------+--------+----------+---------+----------+
|InvoiceNo|StockCode|         Description|Quantity|   InvoiceDate|UnitPrice|    CustomerID|   Country|NumQuant|InvoiceNum|UnitPriceN|CustomerIDN|
+---------+---------+--------------------+--------+--------------+---------+--------------+--------+----------+---------+----------+
|   536401|    21169|YOU'RE CONFUSING ...|       2|12/1/2010 11:21|     1.69|         15862|United Kingdom|       2|    536401|      1.69|     15862|
|   536406|   84406B|CREAM CUPID HEART...|       8|12/1/2010 11:33|     2.75|         17850|United Kingdom|       8|    536406|      2.75|     17850|
|   536408|    20685|DOORMAT RED RETRO...|       2|12/1/2010 11:41|     7.95|         14307|United Kingdom|       2|    536408|      7.95|     14307|
|   536446|    22144|CHRISTMAS CRAFT L...|       2|12/1/2010 12:15|      2.1|         15983|United Kingdom|       2|    536446|      2.10|     15983|
|   536464|    22810|SET OF 6 T-LIGHTS...|       1|12/1/2010 12:23|     2.95|         17968|United Kingdom|       1|    536464|      2.95|     17968|
|   536488|    22468|BABUSHKA LIGHTS S...|       1|12/1/2010 12:31|     6.75|         17897|United Kingdom|       1|    536488|      6.75|     17897|
|   536514|    22866|HAND WARMER SCOTT...|      36|12/1/2010 12:40|      2.1|         17951|United Kingdom|      36|    536514|      2.10|     17951|
|   536529|    21743|STAR PORTABLE TAB...|       6|12/1/2010 13:20|     2.95|         14237|United Kingdom|       6|    536529|      2.95|     14237|
|   536530|    22699|ROSES REGENCY TEA...|       4|12/1/2010 13:21|     2.95|         17905|United Kingdom|       4|    536530|      2.95|     17905|
|   536557|    22114|HOT WATER BOTTLE ...|       2|12/1/2010 14:41|     3.95|         17841|United Kingdom|       2|    536557|      3.95|     17841|
|   536557|    22795|SWEETHEART RECIPE...|       2|12/1/2010 14:41|     6.75|         17841|United Kingdom|       2|    536557|      6.75|     17841|
|   536561|    22866|HAND WARMER SCOTT...|      12|12/1/2010 15:06|      2.1|         12921|United Kingdom|      12|    536561|      2.10|     12921|
|   536630|    21071|VINTAGE BILLBOARD...|       6|12/2/2010 10:56|     1.06|         17850|United Kingdom|       6|    536630|      1.06|     17850|
|   536663|    22737|RIBBON REEL CHRIS...|      20|12/2/2010 12:07|     1.65|         16546|United Kingdom|      20|    536663|      1.65|     16546|
|   536707|    22626|BLACK KITCHEN SCALES|       2|12/2/2010 12:33|      8.5|         12915|United Kingdom|       2|    536707|      8.50|     12915|
|   536741|   85049E|SCANDINAVIAN REDS...|      12|12/2/2010 13:11|     1.25|         13117|United Kingdom|      12|    536741|      1.25|     13117|
|   536762|    22961|JAM MAKING SET PR...|      12|12/2/2010 14:40|     1.45|         16186|United Kingdom|      12|    536762|      1.45|     16186|
|   536762|    22865|HAND WARMER OWL D...|      12|12/2/2010 14:40|      2.1|         16186|United Kingdom|      12|    536762|      2.10|     16186|
|   536804|   84029E|RED WOOLLY HOTTIE...|      48|12/2/2010 16:34|     3.39|         14031|United Kingdom|      48|    536804|      3.39|     14031|
|   536836|    22193|RED DINER WALL CLOCK|       1|12/2/2010 18:08|      8.5|         18168|United Kingdom|       1|    536836|      8.50|     18168|
+---------+---------+--------------------+--------+--------------+---------+--------------+--------+----------+---------+----------+
only showing top 20 rows

>>>
```
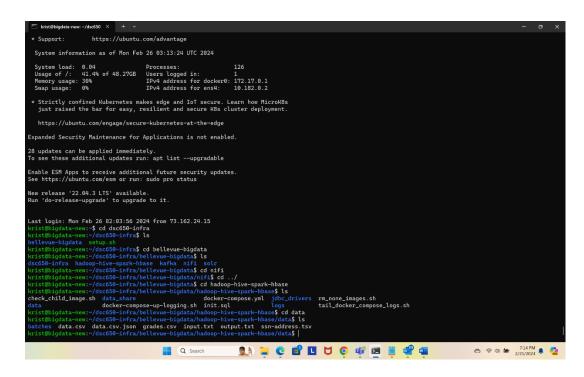
#10:

```
krist@bigdata-new: ~/dsc650
>>> from pyspark.sql.functions import sum, col, desc
>>> df3 = df_any.groupBy("Country") \
...     .agg(sum("NumQuant").alias("sum_quant")) \
...     .sort(desc("sum_quant"))
>>> df3.show(n=100)
+--------------------+---------+
|             Country|sum_quant|
+--------------------+---------+
|      United Kingdom|  4206246|
|         Netherlands|   200361|
|                EIRE|   139320|
|             Germany|   118091|
|              France|   111060|
|           Australia|    83891|
|              Sweden|    36078|
|         Switzerland|    30082|
|               Spain|    27933|
|               Japan|    26016|
|             Belgium|    23237|
|              Norway|    19188|
|            Portugal|    16095|
|             Finland|    10704|
|     Channel Islands|     9485|
|             Denmark|     8235|
|               Italy|     8112|
|              Cyprus|     6340|
|           Singapore|     5241|
|             Austria|     4881|
|              Israel|     3995|
|              Poland|     3684|
|              Canada|     2738|
|             Iceland|     2458|
|                 USA|     2458|
|              Greece|     1557|
|United Arab Emirates|      982|
|               Malta|      970|
|      Czech Republic|      671|
|           Lithuania|      652|
|             Lebanon|      386|
|              Brazil|      356|
|                 RSA|      351|
|             Bahrain|      260|
|        Saudi Arabia|       80|
+--------------------+---------+
```

#11:

```
krist@bigdata-new: ~/dsc650
>>> df5 = df_any.groupBy("Country") \
...     .agg(sum("UnitPriceN").alias("sum_price")) \
...     .sort(desc("sum_price"))
>>> df5.show(n=40)
+--------------------+----------+
|             Country| sum_price|
+--------------------+----------+
|      United Kingdom|1037928.25|
|              France|  36788.80|
|             Germany|  33498.86|
|                EIRE|  32134.30|
|           Singapore|  12949.99|
|               Spain|   9492.39|
|            Portugal|   8636.43|
|             Belgium|   7372.85|
|         Switzerland|   6389.27|
|         Netherlands|   6247.73|
|              Norway|   5662.31|
|             Finland|   3628.44|
|           Australia|   3605.75|
|               Italy|   3576.21|
|              Cyprus|   3466.15|
|     Channel Islands|   3388.00|
|             Austria|   1693.90|
|              Sweden|   1693.69|
|              Poland|   1377.21|
|             Denmark|   1195.55|
|              Canada|    910.48|
|              Israel|    898.53|
|              Greece|    663.29|
|               Japan|    657.21|
|               Malta|    545.19|
|             Iceland|    481.21|
|                 USA|    413.30|
|                 RSA|    248.10|
|             Lebanon|    242.44|
|United Arab Emirates|    229.89|
|              Brazil|    142.60|
|           Lithuania|     99.44|
|             Bahrain|     78.95|
|      Czech Republic|     78.27|
|        Saudi Arabia|     21.16|
+--------------------+----------+
```

#12:



#13:

#14:

```
|UnitPriceN|NumQuant|
+----------+--------+
|      2.55|       6|
|      3.39|       6|
|      2.75|       8|
|      3.39|       6|
|      3.39|       6|
|      7.65|       2|
|      4.25|       6|
|      1.85|       6|
|      1.85|       6|
|      1.69|      32|
|      2.10|       6|
|      2.10|       6|
|      3.75|       8|
|      1.65|       6|
|      4.25|       6|
|      4.95|       3|
|      9.95|       2|
|      5.95|       3|
|      5.95|       3|
|      7.95|       4|
+----------+--------+
only showing top 20 rows

>>> data = assembler.transform(df_num)
>>> data = data.select('features', 'NumQuant')
>>> from pyspark.ml.regression import LinearRegression
>>> lin_reg = LinearRegression(featuresCol='features',
...                            labelCol='NumQuant',
...                            predictionCol='pred_score')
>>> fit = lin_reg.fit(data)
1220375 [Thread-4] WARN  org.apache.spark.ml.util.Instrumentation  - [b8356d43] regParam is zero, which might cause numerical instability and overfitting.
1224478 [dag-scheduler-event-loop] WARN  com.github.fommil.netlib.BLAS  - Failed to load implementation from: com.github.fommil.netlib.NativeSystemBLAS
1224481 [dag-scheduler-event-loop] WARN  com.github.fommil.netlib.BLAS  - Failed to load implementation from: com.github.fommil.netlib.NativeRefBLAS
1359610 [Thread-4] WARN  com.github.fommil.netlib.LAPACK  - Failed to load implementation from: com.github.fommil.netlib.NativeSystemLAPACK
1359611 [Thread-4] WARN  com.github.fommil.netlib.LAPACK  - Failed to load implementation from: com.github.fommil.netlib.NativeRefLAPACK
>>> print(fit.intercept, fit.coefficients)
9.56508375884911 [-0.0027833214440470874]
>>> print(fit.summary.pValues)
[0.36330716528679385, 0.0]
>>> print(fit.summary.r2)
1.5250377750630761e-06
>>>
```

References

Basic Statistics (accessed 2024, March). Basic Statistics - Spark 3.5.0 Documentation

(apache.org)


Convert CSV to JSON, Avro, XML using ConvertRecord (Apache NiFi 1.2+) (Lim, A., 2017,

July). My Cloudera. https://community.cloudera.com/t5/Community-Articles/Convert-CSV-to-

JSON-Avro-XML-using-ConvertRecord-Apache-NiFi/ta-p/246607.


E-Commerce Data (accessed 2024, March). Actual transactions from UK retailer.

https://www.kaggle.com/datasets/carrie1/ecommerce-data.


How to Interpret P-Values in Linear Regression (With Example) (accessed 2024, March).

Statology. How to Interpret P-Values in Linear Regression (With Example) - Statology.


How to Perform Linear Regression in PySpark (With Example) (accessed 2024, March).

How to Perform Linear Regression in PySpark (With Example) - Statology.


Key differences between JSON and CSV (accessed 2024, March). JSON Editor Online.

https://jsoneditoronline.org/indepth/compare/json-vs-

csv/#:~:text=A%20note%20about%20the%20data,keys%20in%20a%20JSON%20file.