# STA323 Assignment1 report

SID: 12110821

Name: ZHANG Chi

## Solution for Q1

### (1)

There are 8 coliuns in the dataset, and the missing value is not allowed when the data is read in. There are total 541909 rows.

The schema of the dataset is defined as follows:

```
schema = StructType([
    StructField("InvoiceNo", IntegerType(), False),      # 第三个：是否允许有空值
    StructField("StockCode", IntegerType(), False),
    StructField("Description", StringType(), False),
    StructField("Quantity", IntegerType(), False),
    StructField("InvoiceDate", StringType(), False),
    StructField("UnitPrice", FloatType(), False),
    StructField("CustomerID", IntegerType(), False),
    StructField("Country", StringType(), False)
])
df = spark.read.csv("./data/Q1_data/retail-dataset.csv", header=True,
    schema=schema)
```

The first 5 rows of the dataset are shown below:

```
+---------+---------+--------------------+--------+--------------+---------+----------+--------------+
|InvoiceNo|StockCode|         Description|Quantity|   InvoiceDate|UnitPrice|CustomerID|       Country|
+---------+---------+--------------------+--------+--------------+---------+----------+--------------+
|   536365|     NULL|WHITE HANGING HEA...|       6|12/1/2010 8:26|     2.55|     17850|United Kingdom|
|   536365|    71053| WHITE METAL LANTERN|       6|12/1/2010 8:26|     3.39|     17850|United Kingdom|
|   536365|     NULL|CREAM CUPID HEART...|       8|12/1/2010 8:26|     2.75|     17850|United Kingdom|
|   536365|     NULL|KNITTED UNION FLA...|       6|12/1/2010 8:26|     3.39|     17850|United Kingdom|
|   536365|     NULL|RED WOOLLY HOTTIE...|       6|12/1/2010 8:26|     3.39|     17850|United Kingdom|
+---------+---------+--------------------+--------+--------------+---------+----------+--------------+
only showing top 5 rows
```

After using `filter` to drop the records where their `quantity` or `UnitPrice` is not positive, 363087 rows are left.

```
df_clean = df.dropna().filter((col("Quantity") > 0) & (col("UnitPrice")> 0))
```

```
+--------+---------+------------------+--------+------------+---------+----------+--------------+
|InvoiceNo|StockCode|       Description|Quantity|  InvoiceDate|UnitPrice|CustomerID|       Country|
+--------+---------+------------------+--------+------------+---------+----------+--------------+
|   536365|    71053| WHITE METAL LANTERN|      6|12/1/2010 8:26|     3.39|     17850|United Kingdom|
|   536365|    22752|SET 7 BABUSHKA NE...|      2|12/1/2010 8:26|     7.65|     17850|United Kingdom|
|   536365|    21730|GLASS STAR FROSTE...|      6|12/1/2010 8:26|     4.25|     17850|United Kingdom|
|   536366|    22633|HAND WARMER UNION...|      6|12/1/2010 8:28|     1.85|     17850|United Kingdom|
|   536366|    22632|HAND WARMER RED P...|      6|12/1/2010 8:28|     1.85|     17850|United Kingdom|
+--------+---------+------------------+--------+------------+---------+----------+--------------+
only showing top 5 rows
```

## (2)

The revenue is calculated by multiplying the `Quantity` and `UnitPrice` for each record, and then summing up the results.

```
1  df_clean.select(expr("sum(UnitPrice * Quantity) as total_cost")).show()
```

```
+----------------+
|      total_cost|
+----------------+
|8015349.50373831|
+----------------+
```

## (3)

Before get top 5 customers that spend most, we need to calculate the total cost for each customer. Then we can sort the total cost in descending order.

Here I use `groupby` and agg to calculate the total cost for each customer, and then use `orderBy` to sort the result.

```
1  df_clean.groupBy("CustomerID").agg(expr("sum(UnitPrice * Quantity) as
   total_cost")).select(col("CustomerID"),col("total_cost")).orderBy("total_cost",a
   scending = False).show(5)
```

```
+----------+------------------+
|CustomerID|        total_cost|
+----------+------------------+
|     14646|265106.91930553317|
|     18102| 253922.7600557804|
|     17450| 180847.0303592682|
|     16446|168472.49374997616|
|     14911|125544.34975004196|
+----------+------------------+
only showing top 5 rows
```
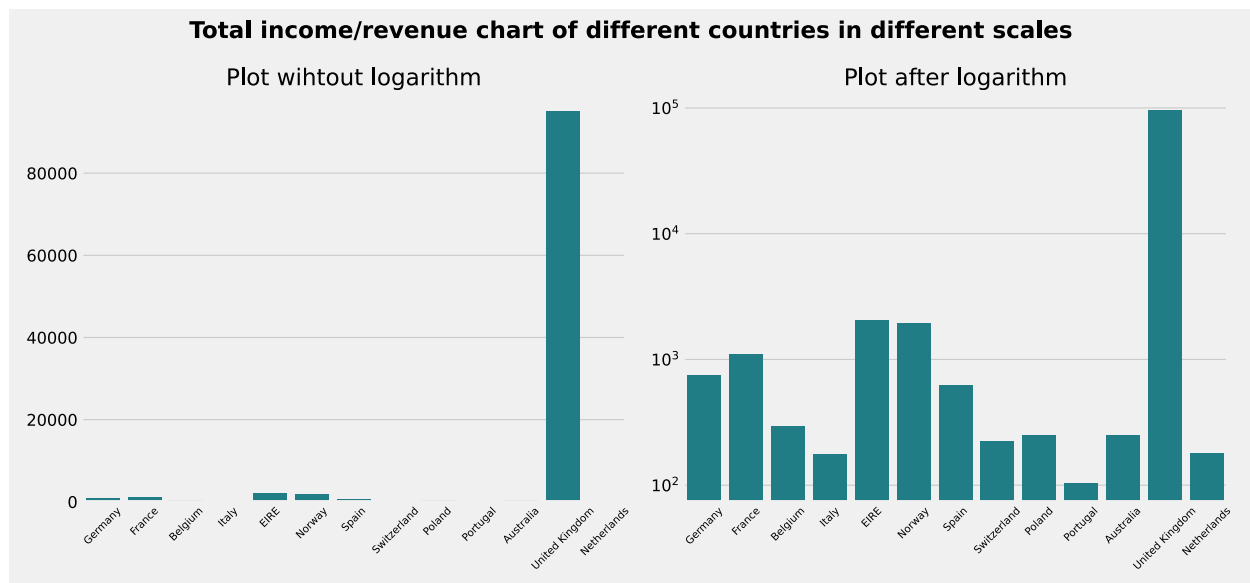
**(4)**

I use `to_timestamp` to convert the `InvoiceDate` to date type. After selecting (by `where` ) the data during designated days, I use `groupBy` and `agg` to calculate the total cost for each day.

```
1  # set the timeParserPolicy to LEGACY
2  spark.conf.set("spark.sql.legacy.timeParserPolicy", "LEGACY")
3  df_4 = df_clean.withColumn("InvoiceDate", to_timestamp(col("InvoiceDate"),
   "MM/dd/yyyy HH:mm"))
4  df_4.where("InvoiceDate between '2010-12-01' and '2010-12-05'
   ").groupBy("Country").agg(expr("sum(UnitPrice * Quantity) as total_cost
   ")).select(col("Country"),col("total_cost")).collect()
```

The bar plot drawn by `seaborn` is shown below. To make it clear, I also draw the plot in `log` scale.



## Solution for Q2

**(1)**

I use `os.listdir` to get the list of all CSV files in the directory. Then I read the first CSV file to initialize the DataFrame, and use `union` to merge the rest of the CSV files. Observing that the first row of each CSV file features the double quote as it beginning, to remove it, the `comment` option is used while loading `.csv` files.

```python
1   import os
2   csv_list = [p for p in os.listdir(path = "data/Q2_data") if p.endswith(".csv")]
3
4   # 读取第一个 CSV 文件来初始化 DataFrame
5   df =
    spark.read.format("csv").option("header","True").option("comment","\"").load(f"
    data/Q2_data/{csv_list[0]}")
6
7   # 从第二个文件开始，读取每一个 CSV 文件并进行合并
8   for i in csv_list[1:]:
9       temp_df =
    spark.read.format("csv").option("header","True").option("comment","\"").load(f"
    data/Q2_data/{i}")
10      df = df.union(temp_df)
```

Before using `spark.read` to load files and remove first lines by `comment`, I tried to use `spark.read.text` and extract the lines by `df.tail()` (except first line) to `createDataFrame`. However, I can't find a easy way to split the each row to make the the number of result correpond to the number of columns as some values in one row is wrapped by `"{}"` with many commas( `,` ) inside. The graphy is shown below. I can't split the row by a single comma.



So that the code below is not used in the final version.
```python
df = spark.read.text("data.csv")
col_name = (df.tail(2)[0]["value"]).split(",")
data = [df.tail(1)[0]["value"].split(",")]
df1 = spark.createDataFrame(data,col_name)
```

Then filter the data to get the result.

```python
1   df.select(col("sequence_alignment_aa"),col("cdr1_aa"),col("cdr2_aa"),col("cdr3_a
    a")).where((length(col("cdr3_aa"))>=10) &  (length(col("cdr3_aa"))<=100 ))
2   df_remove.coalesce(1).write.format("csv").mode("overwrite").save("output/output_
    data/q2_pyspark.csv")
```

```
+-------------------+---------+-------+------------+
|sequence_alignment_aa| cdr1_aa|cdr2_aa|     cdr3_aa|
+-------------------+---------+-------+------------+
| QSVLTQPPSVSGTPGQR...| DSNIGNNF|    KTS| AAWDDPLNAVL|
| GVPDRFSGSTSGTSASL...|     NULL|   NULL|QSFDNSLGGFYV|
| QSALTQPPSASGTPGQR...| RSNIGINT|    SND| DAWDDNLNGPV|
| QSALTQPASMSGSPGQS...|SSDVGASNH|    EVS|  YSYAVGVTFV|
| TLSLTCGSSTGAVHSGY...|TGAVHSGYY|    STD|LLYFGGIQPLWV|
+-------------------+---------+-------+------------+
only showing top 5 rows
```

Save it to a single `.csv` file (Actually, it generatd 4 files in directory `q2_pyspark.csv` ).

```
1   df_remove.coalesce(1).write.format("csv").mode("overwrite").save("output/output_
    data/q2_pyspark.csv")
```

```
📁 q2_pyspark.csv
   📄 _SUCCESS
   📄 ._SUCCESS.crc
   📄 .part-00000-9c1a308a-2505-4f2b-8d74-146ed8211d09-c000.csv.crc
   📄 part-00000-9c1a308a-2505-4f2b-8d74-146ed8211d09-c000.csv
```

## (2)

Run the script by command `bash q2_2.sh` and the output file `q2(2)_sh.csv` is generated. There are 17157 rows (including column names) in the output file `q2(2)_sh.csv`.

```
1  sequence_alignment_aa,cdr1_aa,cdr2_aa,cdr3_aa
2  QSVLTQPPSASGTPGQRVTISCSGSSSNIGSDTVNWFQQLPGSAPKLLIYSNNQRPSGVPDRFSGSKSGTSASLAISGLQSEDEADYYRAAW
   DDSLNGWVVGGGTKLTVL,SSNIGSDT,SNN,AAWDDSLNGWV
3  QSMLTQPPSASGTPGQRVTISCSGGNSNIGSNTVSWYQQFPGAAPKLLIYSSNQRPSGVPARFSGSRSGTSASLAISGLQSEDEAVYYCASW
   DDGLDGFVIFGAGTKLTVL,NSNIGSNT,SSN,ASWDDGLDGFVI
4  QSMLTQPPSASGTPGQRVTISCSGGNSNIGSNTVSWYQQFPGAAPKLLIYSSNQRPSGVPARFSGSRTGTSASLAISGLQSEDEAVYYCASW
   DDGLDGFVIFGAGTKLTVL,NSNIGSNT,SSN,ASWDDGLDGFVI
5  QSMLTQPPSASGTPEQRVTISCSGGNSNIGSNTVSWYQQFPGAAPKLLIYSSNQRPSGVPARFSGSRSGTSASLAISGLQSEDEAVYYCASW
   DDGLDGFVIFGAGTKLTVL,NSNIGSNT,SSN,ASWDDGLDGFVI
6  GVPDRFSGSKSGTSASLAITGLQAEDESAYYCQSYDNSLSVWVFGGGTMLTVL,,,QSYDNSLSVWV
7  QSVLTQPPSASGTPGQRVTISCSGSSSNIGSNYVYWYQQLPGTAPKLLIYRNNQRPSGVPDRFSGSKSGTSASLAISGLRSEDEADYYCAAW
   DDSLSGWVFGGGTKLTVL,SSNIGSNY,RNN,AAWDDSLSGWV
```

## (3)

Run the script by command `bash q2_3.sh` and the output file `q2(3)_sh.csv` is generated. There are 17140 rows (including column names) in the output file `q2(3)_sh.csv`.

```
1  sequence_alignment_aa,cdr1_aa,cdr2_aa,cdr3_aa
2  EIVMTQSPATLSVSPGERATLSCRASQSVSSNLAWYQHKPGQAPRLLIYGTSTRATGIPARFSGSGSGTEFTLTISSLQSEDFA
   VYYCHQYNSWPPGTFGQGTKLEI,QSVSSN,GTS,HQYNSWPPGT
3  DIQMTQSPSSLSASVGDRVTITCRASQSISSYLNWYQQKPGKAPKLLIYAASSLQSGVPSRFSGSGSGTDFTLTISSLQPEDFA
   TYYCQQSYSTHPYTFGQGTKLEI,QSISSY,AAS,QQSYSTHPYT
4  EIVMTQSPATLSVSPGERATLSCRASQSVSSNLAWYQQKPGQAPRLLIYGASTRATGIPARFSGSGSGTEFTLTISSLQSEDFA
   VYYCQQYNNWPPWTFGQGTKVDIK,QSVSSN,GAS,QQYNNWPPWT
5  EIVLAQSPATLSLSPGERATLSCRASQSVSSYLAWYQQKPGQAPRLLIFDASNRATGIPARFSGSGSGTDFTLTISSLEPEDFA
   VYYCQQRNNWPPYTFGQGTKLEI,QSVSSY,DAS,QQRNNWPPYT
6  DIVLTQSPGTLSLSPGERATLSCRATHSINRRFMAWYRQKGGQAPRVIIYGTSIRATGIPDRFSGSGSGTDFTLTISRLEAEDS
   AVYYCQQYDTSQGYPFGQGTKVDIK,HSINRRF,GTS,QQYDTSQGYP
```