

Connect to S3 bucket and read data for MAJOR APPLIANCES REVIEWS in Amazon (data is separated by Tabs compressed as GZ file)

```
% sparksql
%load_in user_data.csv from S3 into a DataFrame
url = "https://s3.amazonaws.com/amazon-reviews-pds/tsv/amazon_reviews_us_Major_Appliances_v1_00.tsv.gz"
spark.sqlContext.addFile(url)

df = spark.read.option("header", "true").csv(SparkFiles.get("amazon_reviews_us_Major_Appliances_v1_00.tsv.gz"), inferSchema=True, sep='\\t', timestampFormat='mm/dd/yyyy')
df.show(10)
```

marketplace	customer_id	review_id	product_id	product_parent	product_title	product_category	star_rating	helpful_votes	total_votes	vine	verified_purchase	review_headline	review_body	review_date
US	16748690	R2EATQZGK3P9J3	B000C5G0X9	63083551R6Z302M_Gallery_	[Major Appliances]	5	1	1	N		Y	If you need a new...[What a great stov...	worked great!	[2015-08-31 00:00:00]
US	15324880	R2EATQZGK3P9J3	B000C5G0X9	81176667Best Hand Clothes...	[Major Appliances]	5	1	0	0	N	Y	Five Stars!		[2015-08-31 00:00:00]
US	16278085	R2EATQZGK3P9J3	B000C5G0X9	34556228Spot SEBIA Ther...	[Major Appliances]	5	0	0	0	N	Y	Fast Shipping!Product exactly what I...		[2015-08-31 00:00:00]
US	15208363	R2EATQZGK3P9J3	B000C5G0X9	561095133Pluriscus JLD003	[Major Appliances]	5	1	1	1	N	Y	Five Stars!Love my refrigera...		[2015-08-31 00:00:00]
US	25414379	R6BZ2DVG60MI	B00017BACU	87426359Javon Bay Portab...	[Major Appliances]	5	0	0	0	N	Y	Five Stars!No more running t...		[2015-08-31 00:00:00]
US	35334175	R1KUCFZFR8E7V0	B000372CM1	29466812Deban Freestand...	[Major Appliances]	5	0	0	0	N	Y	Piece of Junk!It would not cool...		[2015-08-31 00:00:00]
US	48616965	R2EATQZGK3P9J3	B000C5G0X9	183784715Aventi 11c-Volt A...	[Major Appliances]	5	2	2	N		Y	Works awesome for...[Works awesome for...		[2015-08-31 00:00:00]
US	15629126	R1TMDT0RT0P155	B00N0X0V9B	9604763244 Dobby products	[Major Appliances]	5	0	0	0	N	Y	Five Stars!Exactly what I wa...		[2015-08-31 00:00:00]
US	38026019	R2L2UGO2CQCQCQ	B00H7303P2	99247531413 Jet Turb MS&W...	[Major Appliances]	4	0	0	0	N	Y	Four Stars!	As advertised	[2015-08-31 00:00:00]
US	58094294	R1ZBC8CBH8QH	B00N1ESFZ2	1641606The TSSU-08-16 6...	[Major Appliances]	4	0	0	0	N	Y	Ylbut has poor issu...	[It works as advert...	[2015-08-31 00:00:00]

only showing top 10 rows

Interpreter: spark.pyspark. FINISHED Took 34 sec 454 millsec. Updated by Francisco on August 02 2019, 2:28:03 PM (PDT) (outdated)

Count the number of records read from the data source

```
%pyspark
# Rows
df.count()
96981
Interpreter: spark.pyspark. FINISHED Took 1 sec 212 millisec. Updated by Francisco on August 02 2019, 2:28:04 PM (PDT)
```

Count the number of columns available in the data source

```
%pyspark
# Columns
len(df.columns)

15
Interpreter: spark.pyspark. FINISHED Took 159 millise. Updated by Francisco on August 02 2019, 2:28:04 PM (PDT)
```

Select PRODUCT ID and PRODUCT TITLE columns for the products table

```

spark> products = df.select(["product_id", "product_title"])
products.show(5)
-----+-----+
|product_id|product_title|
-----+-----+
|1006796527|G&P38329W Galler...|
|1002556568|Best Hand Cloth...|
|1006452862|Supco SET184 Ther...|
|1008991720|Hideo MS-16081 ...|
|1005778962|Avalon Bay Portab...|
-----+-----+
only showing top 5 rows

```

Interpreter: spark.ySpark. FINISHED Took 174 millisec. Updated by Francisco on August 02 2019, 2:28:05 PM (PDT)

Eliminate duplicated records. A product can be in multiple reviews, we only need 1 record per PRODUCT ID

```
%spark
print(products.count())
products = products.dropDuplicates(["product_id"])
print(products.count())
96901
11694
Interpreter: spark.pyspark. FINISHED Took 2 sec 515 millsec. Updated by Francisco on August 02 2019, 2:28:07 PM (PDT)
```

Identify the number of reviews by CUSTOMER. Group the records by CUSTOMER ID and count the number of records

```

Databricks notebook cell
-- Count the number of customers per country
customers = df.groupby("customer_id").agg({"customer_id": "count"})
customers.show()

-----
(customer_id,count(customer_id))
-----
| 13326061|1|
| 5252251|1|
| 8968913|1|
| 39416583|1|
| 23737233|1|
| 44119972|1|
| 13947800|1|
| 23462558|1|
| 2802853|1|
| 18518845|1|
| 122484|1|
| 14835869|1|
| 20804191|1|
| 16669322|1|
| 24768141|1|
| 24424556|1|
| 44233588|1|
| 13188682|1|
| 23298848|1|
| 58713398|1|
-----
only showing top 20 rows

```

Sort the records grouped by count and rename the column to CUSTOMER_COUNT as defined in the DB schema

```

%sparksql
from pyspark.sql.functions import desc
customers = customers.withColumnRenamed("count(customer_id)", "customer_count")
customers.orderBy(desc("customer_count")).show()

-----+-----
[customer_id][customer_count]
-----+-----
[ 32962888]          161
[ 15480471]          151
[ 5655534]           151
[ 30396577]          141
[ 34553562]          131
[ 29748841]          121
[ 51139148]          121
[ 32230187]          121
[ 51862275]          121
[ 38546694]          111
[ 13649855]          111
[ 40382895]          111
[ 28862753]          111
[ 12914327]          111
[ 48159861]          111
[ 45671123]           91
[ 49484338]           91
[ 24715941]           91
[ 14384334]           91
[ 53874513]           91
-----+-----
only showing top 20 rows

```

Interpreter: spark.pyspark. FINISHED Took 1 sec 312 millisee. Updated by Francisco on August 02 2019, 2:28:10 PM (PDT)

Select columns for table reviews

```

%sparksql
review_id_table = df.select(["review_id", "customer_id", "product_id", "product_parent", "review_date"])
review_id_table.show(5)

+-----+-----+-----+-----+-----+
| review_id|customer_id|product_id|product_parent|review_date|
+-----+-----+-----+-----+-----+
|R1K3Q7H72NAK|161959166|B06074M2ZV|633638551|2015-08-31 00:00:00|
|R1C5L8L3L3|161740808|B000253X86|811766671|2015-08-31 00:00:00|
|R1KX3P4ML7L1A|151220851|B000C45Z86|345642728|2015-08-31 00:00:00|
|R1C2B0W8R9M9C0|320648153|B06WVF1ZG|563657631|2015-08-31 00:00:00|
|R6B120Y7G0U0L1|254144577|B00Y7BM8AC|674236379|2015-08-31 00:00:00|
+-----+-----+-----+-----+-----+
only showing top 5 rows

Interpreter: spark.yespark. FINISHED Took 210 millisec. Updated by Francisco on August 02 2019, 2:28:10 PM (PDT)

```

Extract the DATE from the timestamp column as requested in the table schema

```

%usepark
from pyspark.sql.types import DateType

review_id_table = review_id_table.withColumn("review_date",review_id_table['review_date']).cast(DateType())
#what is the result
review_id_table.show()

.....
| review_id|customer_id|product_id|product_parent|review_date|
.....
|R283WHPF727MAK|15129186|B880783RVCV|633038551|2015-08-31|
|R2AEGUJLEALSP3|16374068|B8802Q5X668|811766671|2015-08-31|
|R1KICD739H8LTLA|15322885|B880F452R6|34562728|2015-08-31|
|R2C2WDFP7PVC0|32884835|B880A1J73G|563052763|2015-08-31|
|R6B1Z0ZY6U001|26414497|B8801Y78MA|874236579|2015-08-31|
|R1KPCZ2NF87E90|35311751|B880332XC1|294467812|2015-08-31|
|R1M9B1E3J0CBRA|38972894|B8805579MA|183784715|2015-08-31|
|R1RT0WYTD0P355|52491265|B880606V8V0|968051524|2015-08-31|
|R1U25U2ZC0C48|48463169|B880173Q01|992475734|2015-08-31|
|R128R5C3J8D9Y0|50794924|B8801S5232|1641680|2015-08-31|
|R1B8NPF3J8RSCZ|3915552|B814513VP1|838188342|2015-08-31|
|R18DNTX0L2P3|17685839|B8802709S8|387104338|2015-08-31|
|R1C8B2N838003|52051058|B88020X6G|188688127|2015-08-31|
|R1B8NPF3J8RSCZ|3915552|B814513VP1|838188342|2015-08-31|
|R2F35P29N90775|17387317|B880LPM008|570132558|2015-08-31|
|R13XNTB8V9Y66|34555366|B88070J5EM|319803847|2015-08-31|
|R1EXUJH5G71TF0|48417244|B880172ALH|364043440|2015-08-31|
|R15974NE8JURV5|19818768|B880C884XK|672722354|2015-08-31|
|R1JCLDH1ET5QLO|43523555|B880F35MG5|137848022|2015-08-31|
|R1V53C2AR24R07|21879631|B8803JG14B|423421857|2015-08-31|
|R23R8V5X239XK|818098|B880LVE9A|386723389|2015-08-31|
.....
only showing top 20 rows

```

Interpreter: spark.pyspark. FINISHED Took 211 millise: Updated by Francisco on August 02 2019, 2:28:10 PM (PDT)

Select columns for the Vine Table

```

Spark
view_table = df.select(["review_id", "star_rating", "helpful_votes", "total_votes", "vine"])
view_table.show(5)

-----
| review_id|star_rating|helpful_votes|total_votes|vine|
-----+-----+-----+-----+----+
|R2830P9E7276AC|5|0|0|N|
|R2EAG1GVLEALSP3|5|1|1|N|
|R3KAC1CP39HLLIA|5|0|0|N|
|R3CZ80F878V0N0|5|1|1|N|
|R6B2Z0ZV6J001|5|0|0|N|
-----
only showing top 5 rows

```

Configure the connection to the DB server hosted in AWS

```
%spark
# Configuration for RDS instance
mode="overwrite"
jdbc_url = "jdbc:postgresql://<Francisco>:5432/my_data_class_db"
config = {"user":"root",
          "password": "*****",
          "driver":"org.postgresql.Driver"}

Interpreter: spark.pyspark. FINISHED Took 109 millise. Updated by Francisco on August 02 2019, 2:28:10 PM (PDT) (outdated)
```

Write the records to the PRODUCTS table

```
%pyspark
# Write DataFrame to table

products.write.jdbc(url=jdbc_url, table='products', mode='append', properties=config)
Interpreter: spark.pyspark. FINISHED Took 27 sec 752 millisecc. Updated by Francisco on August 02 2019, 2:29:12 PM (PDT)
```

Write the records to the REVIEW_ID_TABLE table

```
%pyspark
# Write DataFrame to table

review_id_table.write.jdbc(url=jdbc_url, table='review_id_table', mode=mode, properties=config)
Interpreter: spark.pyspark. FINISHED Took 31 sec 755 millisec. Updated by Francisco on August 02 2019, 2:29:44 PM (PDT)
```

Write the records to the VINE_TABLE table

```
%pyspark
# Write DataFrame to table

products.write.jdbc(url=jdbc_url, table='vine_table', mode=mode, properties=config)
Interpreter: spark.pyspark. FINISHED Took 27 sec 199 millisecond. Updated by Francisco on August 02 2019, 2:30:11 PM (PDT)
```