

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

```
%yspark
from pyspark import SparkFiles
# Load in user_data.csv from S3 into a DataFrame
url = "https://s3.amazonaws.com/amazon-reviews-pds/tsv/amazon_reviews_us_Pet_Products_v1_00.tsv.gz"
spark.sparkContext.addFile(url)

df = spark.read.option("header", "true").csv(SparkFiles.get("amazon_reviews_us_Pet_Products_v1_00.tsv.gz"), inferSchema=True, sep='\t', timestampFormat="mm/dd/yy")

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	28794885	REAC26P879M	B00QK9G04	510387886	(8-Pack) EZhelp ...	Pet Products	5	0	0	N	VIA great purchase ...	Best belly bands ...	2015-08-31 00:00:00
	US	11488981	R3NU7QNZ4HQZEG	B00NKS09W	912374672	Warren Eckstein's...	Pet Products	2	0	1	N	YMy dogs love hugs...	My dogs love hugs...	2015-08-31 00:00:00
	US	43214993	R14QW3XF8QZDP	B00840HJ10	902215727	Tyson's True Chew...	Pet Products	5	0	0	N	YI have been purch...	I have been purch...	2015-08-31 00:00:00
	US	12835065	R2HB7AX0394ZQY	B001G57JK2	568880110	Soft Side Pet Cra...	Pet Products	5	0	0	N	YIt is easy to ope...	It is extremely w...	2015-08-31 00:00:00
	US	26334821	RGMPPQGS4HR3	B004ABH1LG	692846826	EliteField 3-Door...	Pet Products	5	0	0	N	YI dog crate/worked really well...	Dog crate/worked really well...	2015-08-31 00:00:00
	US	22283621	R1D3YVQVCQV6GE	B00AQBLFPM	598674141	Carlson 68-Inch W...	Pet Products	5	0	0	N	YI Five Stars!! I love my gates!	I love my gates!	2015-08-31 00:00:00
	US	14460895	R3V52EAWLFBFQ	B00QFZGZ8	688538603	Dog Seat Cover Mi...	Pet Products	3	0	0	N	YISeat belt tugs on...	Didn't quite work...	2015-08-31 00:00:00
	US	50896354	R3DK083122RQ1	B0001RF9U5	742358789	The Bird Catcher ...	Pet Products	2	0	0	N	YIGreat Pole, but S...	I had the original...	2015-08-31 00:00:00
	US	18449507	R7K4D0KRGKCCG	B0007CB16E	869798483	Cat Bed - Purrfec...	Pet Products	5	1	1	N	N my cat loves it!	The pad is very s...	2015-08-31 00:00:00
	US	50502362	RH1853GAT8Z9F	B00NLS3YZ4	501118658	PetSafe Drinkwell...	Pet Products	5	0	0	N	YI Five Stars! My cat drinks mor...	My cat drinks mor...	2015-08-31 00:00:00

only showing top 10 rows

Interpreter: spark.pyspark. FINISHED Took 1 min 43 sec 76 millsec. Updated by Francisco on August 02 2019, 2:03:50 PM (PDT)

Count the number of records read from the data source

```
%yspark
# Rows
df.count()
```

2643619

Interpreter: spark.pyspark. FINISHED Took 15 sec 233 millsec. Updated by Francisco on August 02 2019, 2:04:05 PM (PDT)

Count the number of columns available in the data source

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

15

Interpreter: spark.pyspark. FINISHED Took 160 millsec. Updated by Francisco on August 02 2019, 2:04:06 PM (PDT)

Select PRODUCT ID and PRODUCT TITLE columns for the products table

```
%yspark
products = df.select(["product_id", "product_title"])
products.show(5)
```

product_id	product_title
B00QK9G04	(8-Pack) EZhelp ...
B00NKS09W	Warren Eckstein's...
B00840HJ10	Tyson's True Chew...
B001G57JK2	Soft Side Pet Cra...
B004ABH1LG	EliteField 3-Door...

only showing top 5 rows

Interpreter: spark.pyspark. FINISHED Took 176 millsec. Updated by Francisco on August 02 2019, 2:04:06 PM (PDT)

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

```
%yspark
print(products.count())
products = products.dropDuplicates(["product_id"])
print(products.count())
```

2643619

239341

Interpreter: spark.pyspark. FINISHED Took 31 sec 918 millsec. Updated by Francisco on August 02 2019, 2:04:38 PM (PDT)

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

```
%yspark
customers = df.groupby("customer_id").agg({"customer_id": "count"})
customers.show()
```

customer_id	count(customer_id)
180270641	1
18765872	1
16711887	1
10742726	2
41169638	1
43622307	1
24540309	2
28253861	1
35329257	2
14552854	1
14525907	5
43502827	5
47282953	1
8201930	1
30309700	2
16485801	4
15966685	21
20840575	2
59048303	1
5596618	1

only showing top 30 rows

Interpreter: spark.pyspark. FINISHED Took 17 sec 337 millsec. Updated by Francisco on August 02 2019, 2:04:55 PM (PDT)

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

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

customer_id	customer_count
2200722	250
33511511	236
16083845	228
25659082	210
23375624	200
43856165	190
5117961	163
3746839	160
16548594	150
52263533	134
50516607	134
33315159	134
18539854	132
18142429	122
29676361	130
39154578	121
14065225	120
38352624	119
42329785	115
50991253	115

only showing top 20 rows

Interpreter: spark.pyspark. FINISHED Took 17 sec 287 millsec. Updated by Francisco on August 02 2019, 2:05:12 PM (PDT)

Select columns for table reviews

```
%yspark
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
REAC26P879M	28794885	B00QK9G04	510387886	2015-08-31 00:00:00
R3NU7QNZ4HQZEG	11488981	B00NKS09W	912374672	2015-08-31 00:00:00
R14QW3XF8QZDP	43214993	B00840HJ10	902215727	2015-08-31 00:00:00
R2HB7AX0394ZQY	12835065	B001G57JK2	568880110	2015-08-31 00:00:00
RGMPPQGS4HR3	26334821	B004ABH1LG	692846826	2015-08-31 00:00:00

only showing top 5 rows

Interpreter: spark.pyspark. FINISHED Took 210 millsec. Updated by Francisco on August 02 2019, 2:05:12 PM (PDT)

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

```
%yspark
from pyspark.sql.types import DateType

review_id_table = review_id_table.withColumn("review_date", review_id_table["review_date"].cast(DateType()))
#Now is the result!
review_id_table.show()
```

review_id	customer_id	product_id	product_parent	review_date
REAC26P879M	28794885	B00QK9G04	510387886	2015-08-31
R3NU7QNZ4HQZEG	11488981	B00NKS09W	912374672	2015-08-31
R14QW3XF8QZDP	43214993	B00840HJ10	902215727	2015-08-31
R2HB7AX0394ZQY	12835065	B001G57JK2	568880110	2015-08-31
RGMPPQGS4HR3	26334821	B004ABH1LG	692846826	2015-08-31
R1D3YVQVCQV6GE	22283621	B00AQBLFPM	598674141	2015-08-31
R3V52EAWLFBFQ	14460895	B00QFZGZ8	688538603	2015-08-31
R3DK083122RQ1	50896354	B0001RF9U5	742358789	2015-08-31
R7K4D0KRGKCCG	18449507	B0007CB16E	869798483	2015-08-31
RH1853GAT8Z9F	50502362	B00NLS3YZ4	501118658	2015-08-31
R33GT7XNP14D4	33938128	B00N0EXMFG	454737777	2015-08-31
R3H7AW0817AYVY	43534294	B001H8QKQY	420905252	2015-08-31
R2ZY4GZKQZK2W6	45555862	B007017H00	302380631	2015-08-31
R2FNHJ3CGM63H0	11147406	B001P3NU30	525778264	2015-08-31
R341G275X0G64	6495078	B002P0H5G5	434117299	2015-08-31
R2H4B0H1S7CFY	2013416	B001P8GSL6	833937831	2015-08-31
R11I0M01N1G293	40459386	B001L8Y598	85343577	2015-08-31
R9B8N0R0H34D	23126000	B011AV430	499241195	2015-08-31
R24H4B0H1PVM	30238471	B0005SH5A1	409532380	2015-08-31
RDC3SH0P955M	35113959	B0003M5C08	259271919	2015-08-31

only showing top 20 rows

Interpreter: spark.pyspark. FINISHED Took 162 millsec. Updated by Francisco on August 02 2019, 2:05:13 PM (PDT)

Select columns for the Vine Table

```
%yspark
vine_table = df.select(["review_id", "star_rating", "helpful_votes", "total_votes", "vine"])
vine_table.show(5)
```

review_id	star_rating	helpful_votes	total_votes	vine
REAC26P879M	5	0	0	N
R3NU7QNZ4HQZEG	2	0	1	N
R14QW3XF8QZDP	5	0	0	N
R2HB7AX0394ZQY	5	0	0	N
RGMPPQGS4HR3	5	0	0	N

only showing top 5 rows

Interpreter: spark.pyspark. FINISHED Took 161 millsec. Updated by Francisco on August 02 2019, 2:05:13 PM (PDT)

Configure the connection to the DB server hosted in AWS

```
%yspark
# 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 millsec. Updated by Francisco on August 02 2019, 2:05:13 PM (PDT)

Write the records to the CUSTOMERS table

```
%yspark
# Write DataFrame to table
customers.write.jdbc(url=jdbc_url, table="customers", mode=mode, properties=config)
```

Interpreter: spark.pyspark. FINISHED Took 2 min 27 sec 832 millsec. Updated by Francisco on August 02 2019, 2:07:41 PM (PDT)

Write the records to the PRODUCTS table

```
%yspark
# Write DataFrame to table
products.write.jdbc(url=jdbc_url, table="products", mode=mode, properties=config)
```

Interpreter: spark.pyspark. FINISHED Took 1 min 5 sec 560 millsec. Updated by Francisco on August 02 2019, 2:08:46 PM (PDT)

Write the records to the REVIEW_ID_TABLE table

```
%yspark
# 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 13 min 38 sec 221 millsec. Updated by Francisco on August 02 2019, 2:22:24 PM (PDT)

Write the records to the VINE_TABLE table

```
%yspark
# Write DataFrame to table
products.write.jdbc(url=jdbc_url, table="vine_table", mode=mode, properties=config)
```

Interpreter: spark.pyspark. FINISHED Took 1 min 4 sec 558 millsec. Updated by Francisco on August 02 2019, 2:23:29 PM (PDT)

