

```

import os
# Find the latest version of spark 3.0 from http://www.apache.org/dist/spark/ and enter as th
# For example:
# spark_version = 'spark-3.0.3'
spark_version = 'spark-3.0.3'
os.environ['SPARK_VERSION']=spark_version

# Install Spark and Java
!apt-get update
!apt-get install openjdk-11-jdk-headless -qq > /dev/null
!wget -q http://www.apache.org/dist/spark/$SPARK_VERSION/$SPARK_VERSION-bin-hadoop2.7.tgz
!tar xf $SPARK_VERSION-bin-hadoop2.7.tgz
!pip install -q findspark

# Set Environment Variables
import os
os.environ["JAVA_HOME"] = "/usr/lib/jvm/java-11-openjdk-amd64"
os.environ["SPARK_HOME"] = f"/content/{spark_version}-bin-hadoop2.7"

# Start a SparkSession
import findspark
findspark.init()

# Download the Postgres driver that will allow Spark to interact with Postgres.
!wget https://jdbc.postgresql.org/download/postgresql-42.2.16.jar

from pyspark.sql import SparkSession
spark = SparkSession.builder.appName("M16-Amazon-Challenge").config("spark.driver.extraClassP

```

▼ Load Amazon Data into Spark DataFrame

```

from pyspark import SparkFiles
url = "https://s3.amazonaws.com/amazon-reviews-pds/tsv/amazon_reviews_us_Furniture_v1_00.tsv"
spark.sparkContext.addFile(url)
df = spark.read.option("encoding", "UTF-8").csv(SparkFiles.get(""), sep="\t", header=True, in
df.show()

```

duct_id	product_parent	product_title	product_category	star_rating	helpful_votes	tc
4HB5E0E	488241329	Shoal Creek Compu...	Furniture	4	0	
42TNMMS	205864445	Dorel Home Produc...	Furniture	5	0	
30MPBZ4	124663823	Bathroom Vanity T...	Furniture	5	1	
5G02ESA	382367578	Sleep Master Ulti...	Furniture	3	0	
5JS8AUA	309497463	1 1/4" GashGuards...	Furniture	3	0	
AVUQQGQ	574537906	Serta Bonded Leat...	Furniture	5	0	
CFY20GQ	407473883	Prepac Shoe Stora...	Furniture	5	2	

FKC48QA	435120460	HomCom PU Leather...	Furniture	5	0
V9IAL9K	356495985	Folding Step Stool	Furniture	5	0
1T4XU1C	243050228	Ace Bayou Adult V...	Furniture	5	0
2HRFLBC	93574483	4D Concepts Audio...	Furniture	5	0
5MISZOC	941823468	Zinus SC-SBBK-14N...	Furniture	5	0
3BMGABC	460567746	Poundex Marble Di...	Furniture	5	1
002VH5Y	829613894	Safavieh Lyndhurs...	Furniture	5	0
LI4RJQ0	816478187	Sauder Boone Moun...	Furniture	5	2
46EC1D0	358594389	Winsome Wood Brea...	Furniture	1	0
00QPL36	312571325	HODEDAH IMPORT Me...	Furniture	3	0
3X7RWB2	402665054	Flash Furniture H...	Furniture	4	0
1TJYPJ8	854989315	Sleep Revolution ...	Furniture	5	0
0TMHX9A	814079288	Flash Furniture V...	Furniture	5	0



▼ Create DataFrames to match tables

```
from pyspark.sql.functions import to_date
# Read in the Review dataset as a DataFrame
```

```
# Create the customers_table DataFrame
customers_df = df.groupby("customer_id").agg({"customer_id": "count"}).withColumnRenamed("count", "customer_count")
customers_df.show()
```

```
+-----+-----+
|customer_id|customer_count|
+-----+-----+
| 17067926|          2|
| 10714827|          1|
| 42560427|          1|
| 30717305|          1|
| 1178966|          1|
| 10429047|          1|
| 41351814|          1|
| 52541790|          2|
| 52512151|          1|
| 37534120|          1|
| 22555935|          1|
| 18681995|          1|
| 2119235|          2|
| 21846356|          1|
| 42251639|          1|
| 7730812|          1|
| 37666248|          1|
| 43676452|          1|
| 41466760|          1|
| 30403003|          1|
+-----+-----+
```

only showing top 20 rows

```
# Create the products_table DataFrame and drop duplicates.
products_df = df.select(['product_id', 'product_title']).drop_duplicates()
products_df.show()
```

```
+-----+-----+
|product_id|      product_title|
+-----+-----+
|B0049H810M|Fun Rugs Surf Tim...|
|B001U0U006|Furniture Repair Set|
|B0076I51JE|Traditional Opera...|
|B007OWPBGU|Frenchi Home Furn...|
|B00DHWGB4M|Circle Design Sid...|
|B00RI6TNJ8|Abrahami Hariz Bu...|
|B00GHZA29Q|nuLOOM Varanas Co...|
|B00GSIIGQS|Milton Greens Sta...|
|B00VZ4RY1I|Dean Shifting San...|
|B007QUM5DM|Charles Petite Sofa|
|B00MN6NTDO|Safavieh Adironda...|
|B00BK31LDQ|Glass Computer De...|
|B0091SXURW|Altra Parsons Des...|
|B00TBVK0Y0|Best Price Mattre...|
|B00A2XM5QC|Legacy Decor Shoj...|
|B002KE7HTQ|Home Styles 5001-...|
|B00V3LVD20|Roundhill Furnitu...|
|B00PN9YSAG|Baxton Studio Hir...|
|B005VAFFN6|Duro Hanley Silve...|
|B001BX1JSC|Flash Furniture V...|
+-----+-----+
```

only showing top 20 rows

```
# Create the review_id_table DataFrame.
# Convert the 'review_date' column to a date datatype with to_date("review_date", 'yyyy-MM-dd
review_id_df = df.select(['review_id', 'customer_id', 'product_id', 'product_parent', to_dat
review_id_df.show()
```

```
+-----+-----+-----+-----+-----+
|      review_id|customer_id|product_id|product_parent|review_date|
+-----+-----+-----+-----+-----+
|R3VR960AHLFKDV|  24509695|B004HB5E0E|  488241329| 2015-08-31|
|R16LGVMFKIUT0G|  34731776|B0042TNMMS|  205864445| 2015-08-31|
|R1AIMEEPYHMOE4|   1272331|B0030MPBZ4|  124663823| 2015-08-31|
|R1892CCSZWZ9SR|  45284262|B005G02ESA|  382367578| 2015-08-31|
|R285P679YWVKD1|  30003523|B005JS8AUA|  309497463| 2015-08-31|
|RLB33HJBXHZHU|  18311821|B00AVUQQGQ|  574537906| 2015-08-31|
|R1VGTZ94DBAD6A|  42943632|B00CFY20GQ|  407473883| 2015-08-31|
|R168KF82ICSOHD|  43157304|B00FKC48QA|  435120460| 2015-08-31|
|R20DIYIJ0OCMOG|  51918480|B00N9IAL9K|  356495985| 2015-08-31|
|RD46RNV0HNZSC|  14522766|B001T4XU1C|  243050228| 2015-08-31|
|R2JD0CETTM3AXS|  43054112|B002HRFLBC|   93574483| 2015-08-31|
```

R33YMW36IDZ6LE	26622950	B006MISZOC	941823468	2015-08-31
R30ZGGUHZ04C1S	17988940	B008BMGABC	460567746	2015-08-31
RS2EZU76IK2BT	18444952	B00C02VH5Y	829613894	2015-08-31
R1GJC1BP028X09	16937084	B00LI4RJQ0	816478187	2015-08-31
R2VKJPGXXEK5GP	23665632	B0046EC1D0	358594389	2015-08-31
R17KS83G3KLT97	4110125	B00DQQPL36	312571325	2015-08-31
R3PQL8SR4NEHWL	107621	B003X7RWB2	402665054	2015-08-31
R2F5WW7WNO5RRG	2415090	B001TJYPJ8	854989315	2015-08-31
R3UDJKVWQCFIC9	48285966	B000TMHX9A	814079288	2015-08-31

+-----+-----+-----+-----+-----+

only showing top 20 rows

```
# Create the vine_table. DataFrame
```

```
vine_df = df.select(['review_id', 'star_rating', 'helpful_votes', 'total_votes', 'vine', 'verified_purchase'])
vine_df.show()
```

review_id	star_rating	helpful_votes	total_votes	vine	verified_purchase
R3VR960AHLFKDV	4	0	0	N	Y
R16LGVFMFKIUT0G	5	0	0	N	Y
R1AIMEEPYHMOE4	5	1	1	N	Y
R1892CCSZWZ9SR	3	0	0	N	Y
R285P679YWVKD1	3	0	0	N	N
RLB33HJBXHZHU	5	0	0	N	Y
R1VGTZ94DBAD6A	5	2	2	N	Y
R168KF82ICSOHD	5	0	0	N	Y
R20DIYIJ00CMOG	5	0	0	N	Y
RD46RNV0HNZSC	5	0	0	N	Y
R2JDOCETTM3AXS	5	0	0	N	Y
R33YMW36IDZ6LE	5	0	0	N	Y
R30ZGGUHZ04C1S	5	1	1	N	Y
RS2EZU76IK2BT	5	0	0	N	Y
R1GJC1BP028X09	5	2	3	N	Y
R2VKJPGXXEK5GP	1	0	0	N	Y
R17KS83G3KLT97	3	0	0	N	Y
R3PQL8SR4NEHWL	4	0	0	N	Y
R2F5WW7WNO5RRG	5	0	0	N	Y
R3UDJKVWQCFIC9	5	0	0	N	Y

+-----+-----+-----+-----+-----+

only showing top 20 rows

▼ Connect to the AWS RDS instance and write each DataFrame to its table.

```
# Store environmental variable
from getpass import getpass
password = getpass('<password>')
# Configure settings for RDS
mode = "append"
jdbc_url="jdbc:<connector ID>:5432/<DB name>"
```

```
config = {"user": "postgres",  
          "password": password,  
          "driver": "org.postgresql.Driver"}
```

📌 password.....

```
# Write review_id_df to table in RDS  
# Took about 3 minutes  
review_id_df.write.jdbc(url=jdbc_url, table='review_id_table', mode=mode, properties=config)  
  
# Write products_df to table in RDS  
# Took about 2 min  
products_df.write.jdbc(url=jdbc_url, table='products_table', mode=mode, properties=config)  
  
# Write customers_df to table in RDS  
# Took 3 minutes 20 seconds  
customers_df.write.jdbc(url=jdbc_url, table='customers_table', mode=mode, properties=config)  
  
# Write vine_df to table in RDS  
# Took 4 minutes 11 seconds  
vine_df.write.jdbc(url=jdbc_url, table='vine_table', mode=mode, properties=config)
```

✓ 3m 47s completed at 2:37 AM

● ✕