Shweta Ajay Shinde
Masters in Data Analytics, San Jose State University
Data 226: Data warehousing
Instructor: Keeyong Han
15th November 2024

Introduction to PySpark

1. Install pyspark and download the 7 log files provided in the demo colab code (+1 pt)

Install pyspark:-

```
| Ipip install pyspark==3.5.3 | Requirement already satisfied: pyspark==3.5.3 in /usr/local/lib/python3.10/dist-packages (3.5.3) | Requirement already satisfied: py4j==0.10.9.7 in /usr/local/lib/python3.10/dist-packages (from pyspark==3.5.3) (0.10.9.7)
```

Download 7 log file:-

```
|wget -nc https://raw.githubusercontent.com/keeyong/sjsu-data226/refs/heads/main/week12/data/sample_web_log_1.log.gz
     lwget -nc https://raw.githubusercontent.com/keeyong/sjsu-data226/refs/heads/main/week12/data/sample_web_log_2.log.gz
     lwget - nc \ https://raw.githubusercontent.com/keeyong/sjsu-data226/refs/heads/main/week12/data/sample\_web\_log\_3.log.gz
     | wget -nc https://raw.githubusercontent.com/keeyong/sjsu-data226/refs/heads/main/week12/data/sample_web_log_4.log.gz
     lwget-nc \\  \  \, https://raw.githubusercontent.com/keeyong/sjsu-data226/refs/heads/main/week12/data/sample\_web\_log\_5.log\_gz
     lwget -nc https://raw.githubusercontent.com/keeyong/sjsu-data226/refs/heads/main/week12/data/sample_web_log_6.log.gz
     2024-11-12 05:46:59 (91.5 MB/s) - 'sample_web_log_2.log.gz' saved [10277610/10277610]
     --2024-11-12 05:46:59-- <a href="https://raw.githubusercontent.com/keeyong/sjsu-data226/refs/heads/main/week12/data/sample_web_log_3.log.gz">https://raw.githubusercontent.com/keeyong/sjsu-data226/refs/heads/main/week12/data/sample_web_log_3.log.gz</a> Resolving raw.githubusercontent.com (raw.githubusercontent.com)... 185.199.110.133, 185.199.108.133, 185.199.109.133, ... Connecting to raw.githubusercontent.com (raw.githubusercontent.com)|185.199.110.133|:443... connected.
     HTTP request sent, awaiting response... 200 OK
     Length: 10276732 (9.8M) [application/octet-stream]
     Saving to: 'sample_web_log_3.log.gz'
     2024-11-12 05:46:59 (90.6 MB/s) - 'sample_web_log_3.log.gz' saved [10276732/10276732]
     --2024-11-12 05:47:00-- https://raw.githubusercontent.com/keeyong/sjsu-data226/refs/heads/main/week12/data/sample_web_log_4.log.gz
     Resolving raw.githubusercontent.com (raw.githubusercontent.com)... 185.199.110.133, 185.199.111.133, 185.199.108.133, ...
     Connecting to raw.githubusercontent.com (raw.githubusercontent.com)|185.199.110.133|:443... connected.
     HTTP request sent, awaiting response... 200 OK
Length: 10277331 (9.8M) [application/octet-stream]
     Saving to: 'sample_web_log_4.log.gz'
     sample_web_log_4.lo 100%[=========>] 9.80M --.-KB/s
     2024-11-12 05:47:00 (95.1 MB/s) - 'sample web log 4.log.gz' saved [10277331/10277331]
```

✓ 2s completed at 9:47 PM

only showing top 20 rows

2. Configure snowflake jar file, set up SparkSession, create an input dataframe (df) and a parsed dataframe (log_df) provided in the demo colab code (+1 pt)

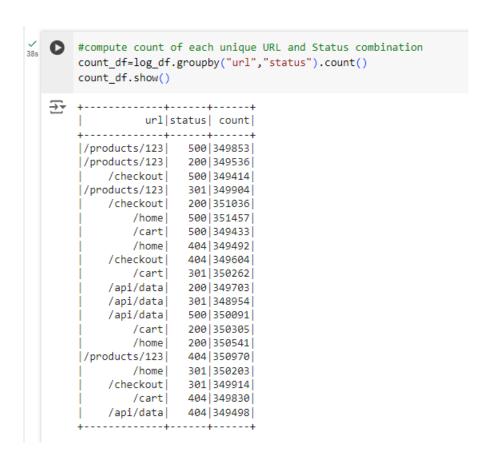
Configure snowflake jar file and set up SparkSession:-

```
[scd /usr/local/lib/python3.10/dist-packages/pyspark/jars && wget https://repo1.maven.org/maven2/net/snowflake/snowflake-jdbc/3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0/snowflake-jdbc-3.19.0
 --2024-11-12 05:55:35-- https://repo1.maven.org/maven2/net/snowflake/snowflake-jdbc/3.19.0/snowflake-jdbc-3.19.0.jar
      Resolving repo1.mayen.org (repo1.mayen.org)... 199.232.192.209, 199.232.196.209, 2a04:4e42:4c::209, ...
       Connecting to repo1.maven.org (repo1.maven.org)|199.232.192.209|:443... connected.
      HTTP request sent, awaiting response... 200 OK
Length: 70986770 (68M) [application/java-archive]
      Saving to: 'snowflake-jdbc-3.19.0.jar'
       2024-11-12 05:55:36 (197 MB/s) - 'snowflake-idbc-3.19.0.jar' saved [70986770/70986770]
                       from pyspark.sql import SparkSession
               [5]
                          import pyspark.sql.functions as F
                          spark = SparkSession.builder.appName("HandleLogFiles").getOrCreate()
                        # Load all .gz files in the directory into a DataFrame
                         df = spark.read.text("*.gz")
               # Check the number of partitions
                 print(df.rdd.getNumPartitions())
                 df.show(truncate=False)
         → 3
                 |value
                  |123.45.67.89 - - [05/Nov/2024:02:08:16 +0000] "DELETE /cart HTTP/1.1" 500 242
                  |192.168.1.1 - - [04/Nov/2024:21:23:39 +0000] "POST /checkout HTTP/1.1" 404 2781
                  |234.56.78.90 - - [05/Nov/2024:07:06:19 +0000] "GET /api/data HTTP/1.1" 301 3758
                  |192.168.1.1 - - [04/Nov/2024:20:03:56 +0000] "POST /home HTTP/1.1" 200 1837
                  |192.168.1.1 - - [04/Nov/2024:21:25:05 +0000] "GET /products/123 HTTP/1.1" 200 3430
                  234.56.78.90 - - [04/Nov/2024:07:38:10 +0000] "GET /api/data HTTP/1.1" 404 3729
                  [123.45.67.89 - - [04/Nov/2024:12:33:22 +0000] "PUT /api/data HTTP/1.1" 404 799
                  |192.168.1.1 - - [04/Nov/2024:07:37:46 +0000] "GET /api/data HTTP/1.1" 500 309
                  123.45.67.89 - - [04/Nov/2024:21:52:36 +0000] "POST /checkout HTTP/1.1" 301 2375
                  | 123.45.67.89 - - [04/Nov/2024:08:36:44 +0000] "DELETE /api/data HTTP/1.1" 404 3449 | 192.168.1.1 - - [05/Nov/2024:03:15:43 +0000] "GET /api/data HTTP/1.1" 200 2319
                  234.56.78.90 - - [05/Nov/2024:01:26:03 +0000] "DELETE /home HTTP/1.1" 500 1168
                  234.56.78.90 - - [05/Nov/2024:03:26:33 +0000] "DELETE /cart HTTP/1.1" 500 1262
                  |123.45.67.89 - - [04/Nov/2024:20:46:25 +0000] "PUT /home HTTP/1.1" 301 4401
                  [123.45.67.89 - - [05/Nov/2024:08:07:51 +0000] "GET /api/data HTTP/1.1" 301 3736
                  123.45.67.89 - - [04/Nov/2024:21:01:30 +0000] "DELETE /cart HTTP/1.1" 404 2418
                  123.45.67.89 - - [04/Nov/2024:09:40:29 +0000] "POST /api/data HTTP/1.1" 301 3260
                  |234.56.78.90 - - [04/Nov/2024:09:23:42 +0000] "GET /home HTTP/1.1" 200 1488
                  |192.168.1.1 - - [04/Nov/2024:11:53:57 +0000] "POST /products/123 HTTP/1.1" 200 2627
                 |234.56.78.90 - - [05/Nov/2024:01:26:01 +0000] "PUT /cart HTTP/1.1" 500 4406
                 +----
```

Extract the information in table format:-

```
\frac{\checkmark}{\Omega_{S}} [7] # Extract the necessary information from log data using regular expressions
       pattern = r'(\d+\.\d+\.\d+\.\d+) - - \[(.*?)\] "(.*?) (.*?) HTTP.*" (\d+) (\d+)'
       log_df = df.select(
          F.regexp_extract("value", pattern, 1).alias("ip"),
          F.regexp_extract("value", pattern, 2).alias("timestamp"),
          F.regexp_extract("value", pattern, 3).alias("method"),
          F.regexp_extract("value", pattern, 4).alias("url"),
          F.regexp_extract("value", pattern, 5).alias("status").cast("integer"),
          F.regexp_extract("value", pattern, 6).alias("size").cast("integer")
  log_df.show()
   timestamp|method| url|status|size|
       |123.45.67.89|05/Nov/2024:02:08...|DELETE| /cart| 500| 242|
       | 192.168.1.1|04/Nov/2024:21:23...| POST|
                                              /checkout | 404 2781
       | 234.56.78.90 | 05/Nov/2024:07:06... | GET | /api/data | 301 | 3758 | | 192.168.1.1 | 04/Nov/2024:20:03... | POST | /home | 200 | 1837 |
       | 192.168.1.1|04/Nov/2024:21:25...| GET|/products/123| 200|3430
       234.56.78.90 04/Nov/2024:07:38... | GET | /api/data | 404 3729 |
       | 123.45.67.89 | 04/Nov/2024:12:33... | PUT | /api/data | 404 | 799 |
       | 192.168.1.1|04/Nov/2024:07:37...| GET| /api/data| 500| 309|
       |123.45.67.89|04/Nov/2024:21:52...| POST| /checkout| 301|2375|
       |123.45.67.89|04/Nov/2024:08:36...|DELETE| /api/data| 404|3449|
       | 192.168.1.1|05/Nov/2024:03:15...| GET| /api/data| 200|2319|
                                               /home | 500 | 1168 |
       234.56.78.90 05/Nov/2024:01:26... | DELETE
       234.56.78.90 05/Nov/2024:03:26... | DELETE |
                                                  /cart| 500|1262
                                                /home| 301|4401|
       |123.45.67.89|04/Nov/2024:20:46...| PUT| | | |
       |123.45.67.89|05/Nov/2024:08:07...| GET| /api/data| 301|3736|
       |123.45.67.89|04/Nov/2024:21:01...|DELETE| /cart| 404|2418|
       |123.45.67.89|04/Nov/2024:09:40...| POST| /api/data| 301|3260|
       234.56.78.90|04/Nov/2024:09:23...| GET| /home| 200|1488|
       | 192.168.1.1|04/Nov/2024:11:53...| POST|/products/123| 200|2627|
       |234.56.78.90|05/Nov/2024:01:26...| PUT| /cart| 500|4406|
       +----+
       only showing top 20 rows
```

3. Compute the counts of url and status combination. Use DataFrame operations to compute the count of each unique url and status combo (+2 pt).

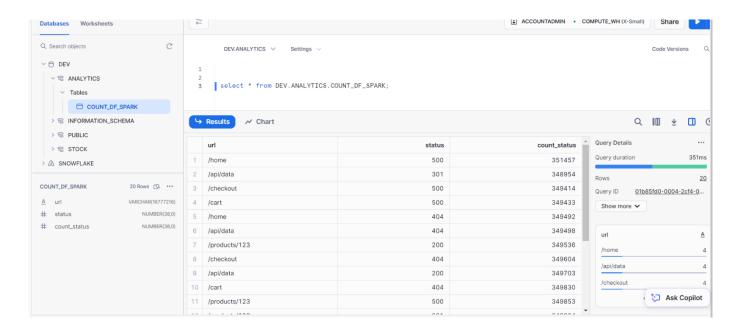


4. Compute the same (step 3) with SparkSQL (+2 pt)

```
# Register the DataFrame as a temporary SQL table
   log_df.createOrReplaceTempView("logs_df_status")
   # Use SparkSQL to count URLs with 404 status
   count_df_spark = spark.sql("""
       SELECT url, status, COUNT(*) as count status
       FROM logs_df_status
       GROUP BY url, status
       ORDER BY count status DESC
   count_df_spark.show()
   +----+
           url|status|count_status|
   +----+
          /home| 500| 351457|
       /checkout| 200|
                          351036
    |/products/123| 404|
                          350970
           /home| 200|
                         350541
           /cart| 200|
                         350305
           /cart| 301|
                         350262
                         350203
           /home| 301|
                          350091
        /api/data| 500|
                         349914
        /checkout| 301|
                         349904
    |/products/123| 301|
                         349853
    |/products/123| 500|
                         349830
           /cart| 404|
        /api/data| 200|
                          349703
        /checkout| 404|
                          349604
                 200 |
404 |
    /products/123|
                          349536
        /api/data|
                           349498
           /home
                  404
                           349492
           /cart|
                  500
                           349433
                 500
                          349414
        /checkout
       /api/data|
                 301
                          348954
```

5. Load this dataframe from step 4 into a table in Snowflake (+2 pt)

6. Capture the screenshot of the table (SELECT results) in Snowflake (+1 pt)



Git Link:

 $\underline{https://github.com/ShindeShwetaK/DW_Assignment10/blob/main/Shweta_Shinde_Hom_ework10.ipynb}$