# Assignment-10

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



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)

## Create a parsed dataframe (log_df)

```python
# 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")
)
```

```
[9] log_df.show()
```

```
+-----------+--------------------+------+-------------+------+----+
|         ip|           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|
```

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

## Let's compute top 404 urls

```python
# Keep only 404 error logs
error_404_logs = log_df.filter(log_df.status == 404)
```

```python
# Group by URL and then count, and sort by count in descending order
url_404_count = error_404_logs.groupBy("url").count().orderBy(F.desc("count"))
```

```
[12] # print the outcome
     url_404_count.show()
```

```
+-------------+------+
|          url| count|
+-------------+------+
|/products/123|350970|
|        /cart|349830|
|    /checkout|349604|
|    /api/data|349498|
|        /home|349492|
+-------------+------+
```

Compute the same (step 3) with SparkSQL (+2 pt)
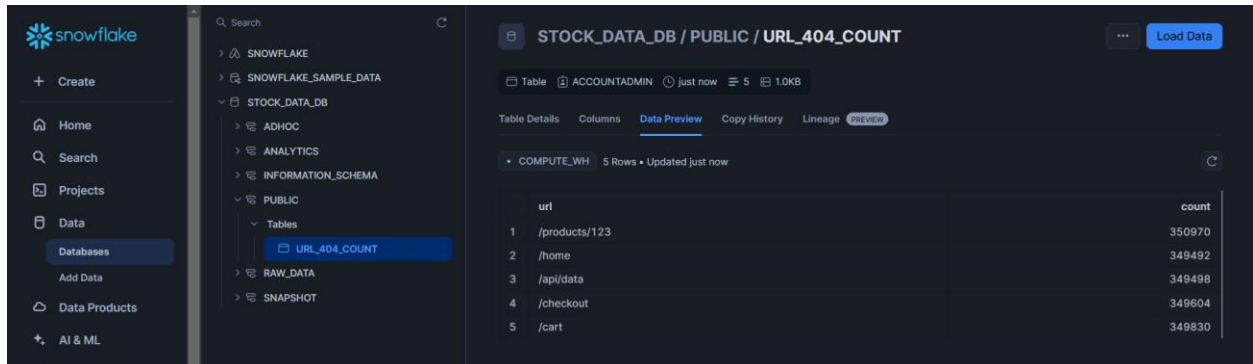
## Now Let's do this in SparkSQL

```
[13] # Register the DataFrame as a temporary SQL table
     log_df.createOrReplaceTempView("logs")
```

```python
# Use SparkSQL to count URLs with 404 status
url_404_count = spark.sql("""
    SELECT url, COUNT(1) as count
    FROM logs
    WHERE status = 404
    GROUP BY url
    ORDER BY count DESC
""")
```
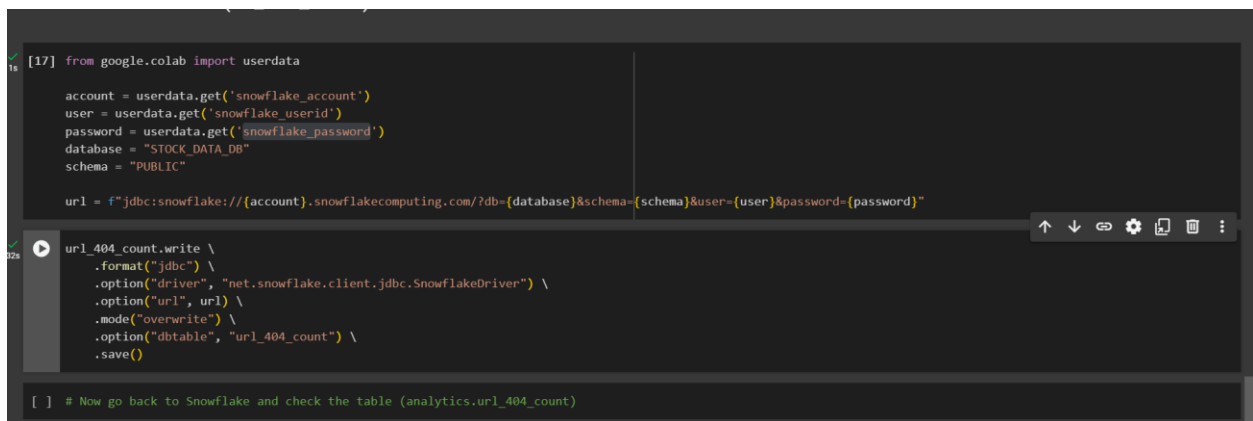
```
[15] url_404_count.show()
```

```
+-------------+------+
|          url| count|
+-------------+------+
|/products/123|350970|
|        /cart|349830|
|    /checkout|349604|
|    /api/data|349498|
|        /home|349492|
```

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





https://colab.research.google.com/drive/1rYIHzsvERLNVe3SUv9fbv1_bHK5d6aL9?usp=sharing