

2021538

Romaisa

Assignment 03

CE408

Downloading the Spark Container:

```
C:\Windows\system32\cmd.exe
D:\Semester07\CE408\A03\spark-netflix-eda>docker run -it --rm apache/spark-py
docker: error during connect: Head "http://%2F%2F.%2Fpipe%2FdockerDesktopLinuxEngine/_ping": open //./pipe/dockerDesktopLinuxEngine: The system cannot find the file specified.
See 'docker run --help'.

74ac377868f8: Downloading [=====>] 10.49MB/30.43MB
a182a611d05b: Downloading [====>] 1.049MB/16.98MB
a182a611d05b: Downloading [=====>] 8.389MB/16.98MB
af739662ba28: Downloading [==>] 6.291MB/121.5MB
af739662ba28: Downloading [====>] 6.291MB/121.5MB
74ac377868f8: Downloading [=====>] 10.49MB/30.43MB
c0b51f9f6377: Downloading [==>] 2.097MB/46.94MB
4f4fb00ef54: Already exists
78ff09ab9807: Download complete
74ac377868f8: Download complete
af739662ba28: Downloading [====>] 6.291MB/121.5MB
a182a611d05b: Download complete
af739662ba28: Download complete
1365af4d8cae0: Download complete
a09fe30821a3: Downloading [=====>] 1.049MB/1.75MB
c0b51f9f6377: Download complete
b0b2797952c6: Download complete
54c642dfa726: Downloading [=====>] 3.146MB/21.33MB
a09fe30821a3: Downloading [=====>] 1.049MB/1.75MB
b235cf4705af: Pulling fs layer
a09fe30821a3: Download complete
54c642dfa726: Download complete
af739662ba28: Downloading [====>] 1.049MB/121.5MB
a182a611d05b: Downloading [=====>] 2.097MB/16.98MB
ac93fae56125: Download complete
e072aad951d: Downloading [=====>] 1.049MB/1.614MB

Digest: sha256:bef1ed7818dd775c8a88224d5b2550c9a85ff81860f76b44e5357abdd849bb5 2.097MB/16.98MB
Status: Downloaded newer image for apache/spark-py:latest
cd /app1d05b: Downloading [=====>] 2.097MB/16.98MB
++ id -u05b: Downloading [=====>] 2.097MB/16.98MB
++ myuid=185b: Downloading [=====>] 2.097MB/16.98MB
++ id -g
++ mygid=0
+ set +e
++ getent passwd 185ading [>] 1.049MB/307.7MB
+ uidentry=
+ set -e
+ '[' -z '' ']'
+ '[' -w /etc/passwd ']'
+ echo '185x:185:0:anonymous uid:/opt/spark:/bin/false'
+ '[' -z /opt/java/openjdk ']'
+ SPARK_CLASSPATH="/opt/spark/jars/*"
+ envad951d: Download complete
+ grep SPARK_JAVA_OPT_
+ sed 's/[~]*~\(\.\.\)/\1/g'
+ sort -t_ -k4 -n
++ command -v readarray
+ '[' readarray ']'
+ readarray -t SPARK_EXECUTOR_JAVA_OPTS
+ '[' -n '' ']'
+ '[' -z ']'
+ '[' -z ']'
+ '[' -n '' ']'
+ '[' -z ']'
+ '[' -z ']'
+ '[' -z ']'
+ '[' -z ']'
+ SPARK_CLASSPATH="/opt/spark/conf:/opt/spark/jars/*"
+ case "$1" in
+ echo 'Non-spark-on-k8s command provided, proceeding in pass-through mode...'


```

Running the Script in Spark

```
(base) PS C:\Users\hp> docker run -it --v /d/Semester07/CE408/A03/spark-netflix-eda:/data apache/spark-py /opt/spark/bin/pyspark
++ id -g
++ myuid=185
++ id -p
++ mygid=0
+ set +e
++ getent passwd 185
+ uidentry=
+ set -e
+ '[' -z '' ']'
+ '[' -w /etc/passwd ']'
+ echo '185x:185:0:anonymous uid:/opt/spark:/bin/false'
+ '[' -z /opt/java/openjdk ']'
+ SPARK_CLASSPATH="/opt/spark/jars/*"
+ env
+ grep SPARK_JAVA_OPT_
+ sort -t_ -k4 -n
+ sed 's/[~]*~\(\.\.\)/\1/g'
++ command -v readarray
+ '[' readarray ']'
+ readarray -t SPARK_EXECUTOR_JAVA_OPTS
+ '[' -n '' ']'
+ '[' -z ']'
+ '[' -z ']'
+ '[' -n '' ']'
+ '[' -z ']'
+ '[' -z ']'
+ '[' -z ']'
+ SPARK_CLASSPATH="/opt/spark/conf:/opt/spark/jars/*"
+ case "$1" in
+ echo 'Non-spark-on-k8s command provided, proceeding in pass-through mode...'
Non-spark-on-k8s command provided, proceeding in pass-through mode...
+ CMD=""
+ exec /usr/bin/tlnt -s -- /opt/spark/bin/pyspark
Python 3.10.6 (main, Mar 10 2023, 10:55:28) [GCC 11.3.0] on linux

RAM 2.45 GB CPU 1.89% Disk 1018.86 GB avail. of 1081.10 GB
```

»»» ↓

	mean	null	null	1124.7692307692307	null	null	1944.0	null 2014.189598270172	2016.8	1994.0	↓
--	------	------	------	--------------------	------	------	--------	------------------------	--------	--------	---

A diagram showing a horizontal line with 11 vertical tick marks. The top line is solid, and the bottom line is dashed. A downward arrow is on the right.

# Terminal

```
>>> # Show unique values in 'type' column
>>> df.select("type").distinct().show()
+-----+
|      type|
+-----+
|      null|
|   TV Show|
|     Movie|
|William Wyler|
+-----+

>>>
>>> # Count the number of each type (movie or tv show)
>>> df.groupBy("type").count().show()
+-----+-----+
|      type|count|
+-----+-----+
|      null|    1|
|   TV Show| 2676|
|     Movie| 6131|
|William Wyler|    1|
+-----+-----+

>>>
>>> # Analyze distribution of 'rating' values
>>> df.groupBy("rating").count().show()
+-----+-----+
|      rating|count|
+-----+-----+
| November 1, 2028|    1|
| Shavidee Trotter|    1|
| Adriane Lenox|    1|
|      TV-Y|   307|
| Maury Chaykin|    1|
|      2019|    1|
```

