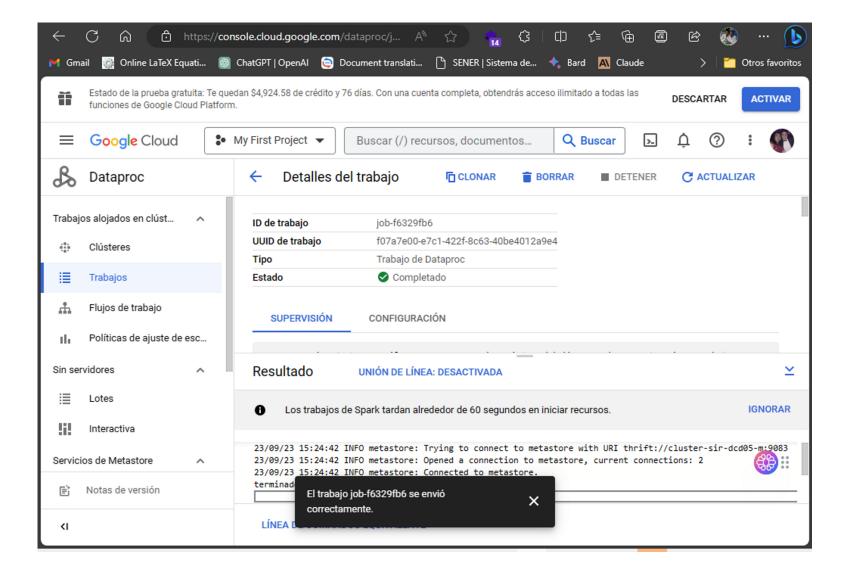# Diplomado en Ciencia de Datos UNAM
# Módulo 13 Datos Masivos
# Sesión 3 del 23 Septiembre de 2023
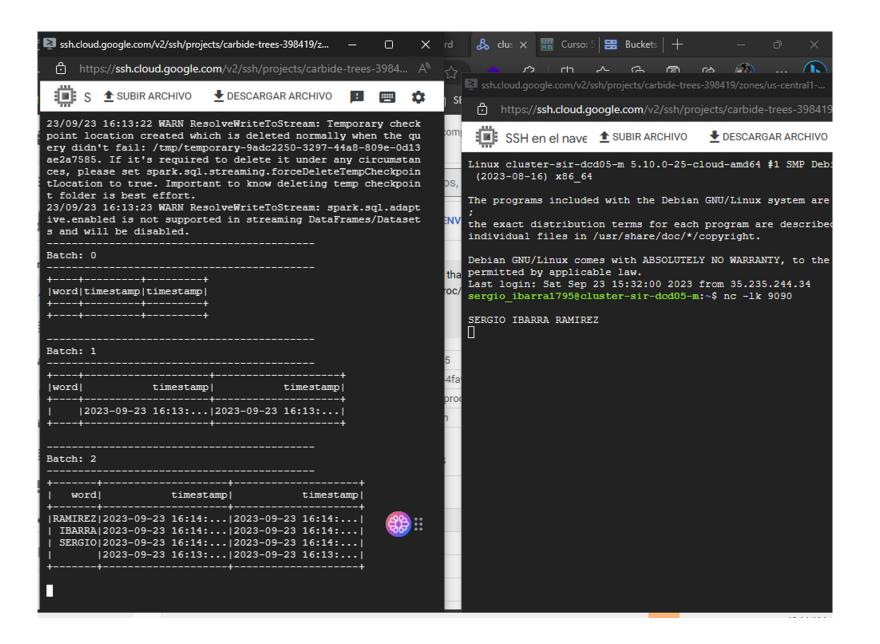
**Sergio Ibarra Ramírez**

# Pantalla de Job de pysparkSQL ejecutado

# Pantalla 2



```
ivysettings.xml file not found in HIVE_HOME or HIVE_CONF_DIR
,/etc/hive/conf.dist/ivysettings.xml will be used
23/09/23 15:36:02 INFO SparkEnv: Registering MapOutputTracke
r
23/09/23 15:36:02 INFO SparkEnv: Registering BlockManagerMas
ter
23/09/23 15:36:02 INFO SparkEnv: Registering BlockManagerMas
terHeartbeat
23/09/23 15:36:02 INFO SparkEnv: Registering OutputCommitCoo
rdinator
Spark master: yarn, Application Id: application_169547971637
8_0002
2020    6       64
2020    11      18
2020    3       36
2020    9       24
2020    12      32
2020    4       44
2021    3       2
2021    2       10
2020    1       19428
2020    8       42
2021    1       34
2020    7       52
2020    2       132
2020    10      40
2020    5       42
Time taken: 12.108 seconds, Fetched 15 row(s)
sergio_ibarra1795@cluster-sir-dcd05-m:~$ hdfs dfs -ls hdfs:/
//tmp/dcd/job/covid
Found 3 items
-rw-r--r--   2 sergio_ibarra1795 hadoop          0 2023-09-2
3 15:33 hdfs:///tmp/dcd/job/covid/_SUCCESS
-rw-r--r--   2 sergio_ibarra1795 hadoop      23886 2023-09-2
3 15:33 hdfs:///tmp/dcd/job/covid/part-00000-3233f9b8-c   -4
fb4-b5fd-89dbf1ecbdea-c000.snappy.parquet
-rw-r--r--   2 root              hadoop      23886 2023-09-2
3 15:24 hdfs:///tmp/dcd/job/covid/part-00000-a629d3cf-ea1f-4
ff6-9f9a-c0bf253c6740-c000.snappy.parquet
sergio_ibarra1795@cluster-sir-dcd05-m:~$
```

# Pantalla 3

# Pantalla 4