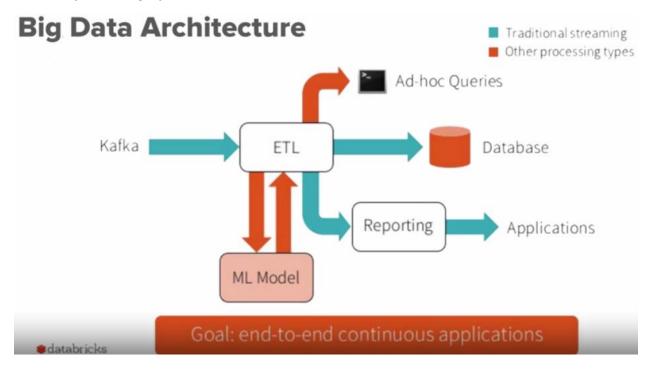
University of California Davis and Databricks collaboration Distributed Computing with Spark SQL

103 Spark SQL Engineering Data Pipelines

Pipelining: Process of moving data through an application

Data Ecosystems: Legacy Data Warehouses, Databases, Data Lakes



Architecture Needs: Scalability

A scalable way to receive data

Producers publish to topics

Consumers consume from topics

Apache Kafka, Amazon Kinesis, Azure Event Hub

Architecture Needs: ETL Process

Extract - Transform - Load

Raw data → usable data

Other Architecture Needs

Transactional database

Ad hoc queries

Machine learning

Reporting

Scalable data lake

Decoupled Storage and Compute

No data stays on the cluster

Resources are elastic

Updates are easy

IO vs. CPU Bound Problems

IO is network transfer (data transfer)

CPU is computation (processing data)

Data transfer is normally the bottleneck for tasks

Spark solves this by colocating your storage

Data Sources

Traditional databases like Postgres, SQL Server, and MySQL

NoSQL databases/documents stores like MongoDB

Distributed databases like Cassandra and Redshift

Blob stores like S3 and the Azure Blob

Message brokers like Kafka and Kinesis

Data warehouses like Hive

File types like CSV, Parquet, and Avro

Online Analytical Processing (OLAP)

Reporting

Ad hoc analysis

Blob stores

Data warehouses

Online Transaction Processing (OLTP)

Databases

Web traffic

Streaming

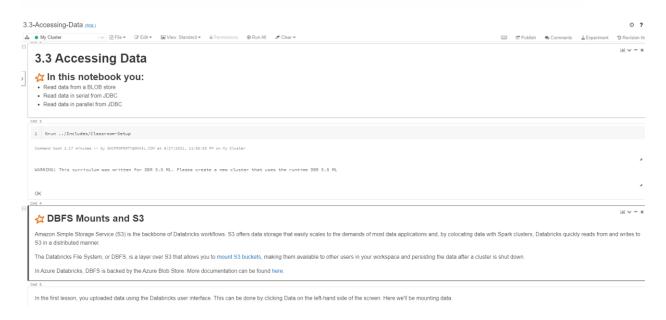
Connecting to two common data sources: BLOB and JDBC

BLOB and JDBC

BLOB is a common way of storing large amounts of data

JDBC is an application programming interface (API) for Java environments

Predicate push down functionality

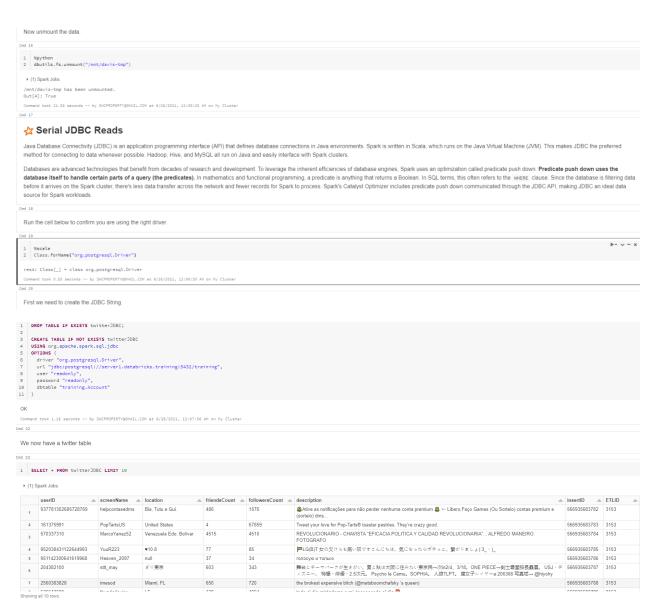


Define your AWS credentials. Below are defined read-only keys, the name of an AWS bucket, and the mount name we will be referring to in DBFS.

For getting AWS keys, take a look at take a look at the AWS documentation



You now have access to an unlimited store of data. This is a read-only S3 bucket. If you create and mount your own, you can write to it as well



± 1 h. m

Command took 5.90 seconds -- by SWCPROPERTY@GMAIL.COM at 6/28/2021, 12:07:41 AM on My Cluster

Add a subquery to dbtable , which pushes the predicate to JDBC to process before transferring the data to your Spark cluster.

7 123891658 itsgirlie1818 Phillippines 524 1094 i love music and i'm a fangirl :)

```
1 DROP TABLE IF EXISTS twitterPhilippinesJDBC;
 3 CREATE TABLE twitterPhilippinesJDBC
4 USING org apachs
       CREATE TABLE twitterPhilippinesJDBC
USING org.apache.spark.sql.jdbc
OPTIONS (
driver "org.postgresql.Driver",
url "jdbc:postgresql://serverl.databricks.training;5432/training",
user "readonly",
password "readonly",
pessions recounty,

dbtable "(SELECT * FROM training.Account WHERE location = 'Philippines') as subq"

11 )
 Command took 0.74 seconds -- by SWCPROPERTY@GMAIL.COM at 6/28/2021, 12:11:04 AM on My Cluster
1 SELECT * FROM twitterPhilippinesJDBC LIMIT 10
   ▶ (1) Spark Jobs
               userID 	riangle screenName 	riangle location 	riangle friendsCount 	riangle followersCount 	riangle description

        userID
        screenName
        location
        friendsCount
        followersCount

        1
        877436554766987264
        m0wskito
        Philippines
        224
        52

        2
        2807579618
        destinedtoV
        Philippines
        158
        72

        3
        422488190
        erickmacasaet
        Philippines
        349
        456

        4
        54405827
        asdfghkjirah
        Philippines
        1837
        954

        5
        945951715689250818
        seokpiggy
        Philippines
        512
        171

        6
        384854316
        lovatic_McDLove
        Philippines
        665
        306

                                                                                                                                                                                                                                                                                                                                                                558345749937 3153
                                                                                                                                                                                                       schatz
                                                                                                                                                                                                       But I fell for you so bad. Though I knew that you and I weren't a good fit (ctto: pic)
                                                                                                                                                                                                                                                                                                                                                                549755819681 3153
                                                                                                                                                                                                      Broke but still Grinding at my 20's
                                                                                                                                                                                                                                                                                                                                                                566935684829 3153
```

Jeremiah 29:11 x A wandering soul that just wants to go to Seoul

Obsessing over Demi Lovato, Girls' Generation, APink, IU, BlackPink 💖

My handsomeness is staining my life.

575525617821 3153 558345751013 3153

575525618268 3153

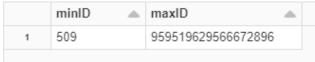
549755821532 3153

Showing all 10 rows. ± h. =

Command took 1.16 seconds -- by SWCPROPERTY@GMAIL.COM at 6/28/2021, 12:11:19 AM on My Cluster







Showing all 1 rows.



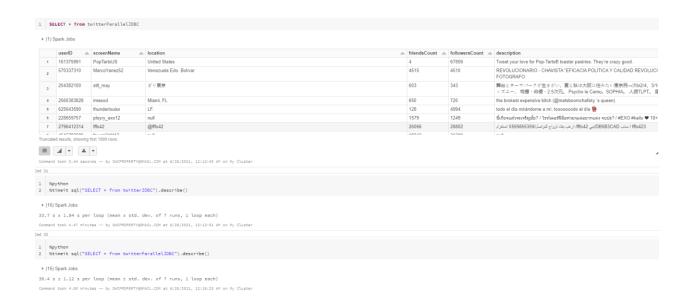
Command took 2.76 seconds -- by SWCPROPERTY@GMAIL.COM at 6/28/2021, 12:12:01 AM on My Cluster

```
Cmd 29
```

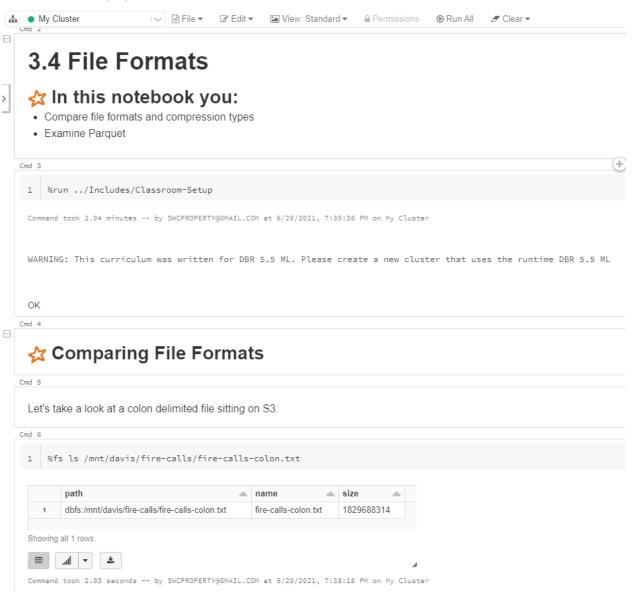
```
DROP TABLE IF EXISTS twitterParallelJDBC;
1
2
   CREATE TABLE IF NOT EXISTS twitterParallelJDBC
3
   USING org.apache.spark.sql.jdbc
4
5
   OPTIONS (
     driver "org.postgresql.Driver",
6
     url "jdbc:postgresql://serverl.databricks.training:5432/training",
7
     user "readonly",
8
9
    password "readonly",
     dbtable "training.Account",
10
      partitionColumn '"userID"',
11
12
     lowerBound 2591,
13
     upperBound 951253910555168768,
      numPartitions 25
14
15
```

ΟK

Command took 0.59 seconds -- by SWCPROPERTY@GMAIL.COM at 6/28/2021, 12:12:07 AM on My Cluster







Take a look at the first few lines of the file.

```
Truncated to first 1808 bytes:

[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes:
[Truncated to first 1808 bytes
```



Showing all 34 rows.



Are these data types correct? All of them are string types.

We need to tell Spark to infer the schema.

```
Cmd 14
     CREATE OR REPLACE TEMPORARY VIEW fireCallsCSV
 1
     USING CSV
 3
     OPTIONS (
 4
          path "/mnt/davis/fire-calls/fire-calls-colon.txt",
          header "true",
 5
 6
         sep ":",
 7
          inferSchema "true"
 8
  ▶ (2) Spark Jobs
 OK
 Command took 1.04 minutes -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 7:41:19 PM on My Cluster
Cmd 15
```

Now take a look at how Spark inferred the data types.

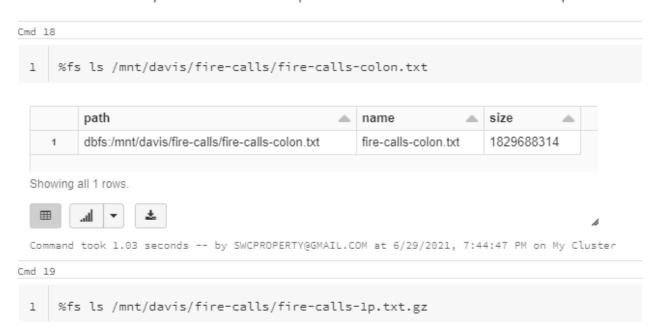
```
1 DESCRIBE fireCallsCSV
```

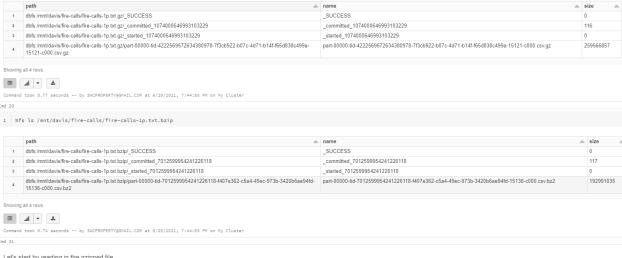


Wow, that took a long time just to figure out the schema for this file!

Now let's try the same thing with compressed files (Gzip and Bzip formats).

Notice that the bzip file is the most compact - we will see if it is the fastest to operate on.





Let's start by reading in the gzipped file.

```
1 CREATE OR REPLACE TEMPORARY VIEW fireCallsCSVgzip
3 OPTIONS (
       path "/mnt/davis/fire-calls/fire-calls-lp.txt.gz",
5
        header "true",
       sep ":",
       inferSchema "true"
8 )
 ▶ (2) Spark Jobs
```

Command took 39.74 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 7:44:59 PM on My Cluster

Wow! That took way longer than inferring the schema on the uncompressed data. Even though it took up less storage space, we had to pay for that in computation.

You'll notice that the resulting view is comprised of only 1 partition, which makes this data very slow to query later on.

```
1 %python
2 sql("SELECT * FROM fireCallsCSVgzip").rdd.getNumPartitions()
Command took 0.24 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 7:48:16 PM on My Cluster
```

Let's compare the speed of reading in the gzip file to the bzip file!

```
Cmd 26
 1 CREATE OR REPLACE TEMPORARY VIEW fireCallsCSVbzip
   USING CSV
 3 OPTIONS (
       path "/mnt/davis/fire-calls/fire-calls-lp.txt.bzip",
       header "true",
       sep ":",
        inferSchema "true"
```

```
1  %python
2  sql("SELECT * FROM fireCallsCSVbzip").rdd.getNumPartitions()

Out[3]: 8
Command took 0.18 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 7:51:09 PM on My Cluster

Cmd 28
```

Bzip is a "splittable" file format, so it is much better to use than gzip when working with row-based formats for querying later on.

Now let's go ahead and compare that to reading in from a columnar format: Parquet.

1 DESCRIBE fireCallsParquet

	col_name	data_type 🔺	comment 🔺
1	Call_Number	int	null
2	Unit_ID	string	null
3	Incident_Number	int	null
4	Call_Type	string	null
5	Call_Date	string	null
6	Watch_Date	string	null
7	Received DtTm	string	null

Showing all 34 rows.

Look at how fast it is to get the schema from a Parquet file! That is because the Parquet file stores the data and the associated metadata. Cmd 32 Compare the performance between the three file types. We are going to use a Python helper function called timeit to calculate how long the query takes to execute. 1 %python parquetDF = sql("SELECT * FROM fireCallsParquet") 3 %timeit -nl -rl parquetDF.select("City").where("City == 'San Francisco'").count() ▶ (2) Spark Jobs ▶ ■ parquetDF: pyspark.sql.dataframe.DataFrame = [Call_Number: integer, Unit_ID: string ... 32 more fields] 7.87 s \pm 0 ns per loop (mean \pm std. dev. of 1 run, 1 loop each) Command took 8.05 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 7:52:05 PM on My Cluster Cmd 34 1 %python csvDF = sql("SELECT * FROM fireCallsCSV") 3 %timeit -n1 -r1 csvDF.select("City").where("City == 'San Francisco'").count() ▶ (2) Spark Jobs ▶ ■ csvDF: pyspark.sql.dataframe.DataFrame = [Call Number: integer, Unit ID: string ... 32 more fields] 40.9 s \pm 0 ns per loop (mean \pm std. dev. of 1 run, 1 loop each) Command took 41.10 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 7:52:07 PM on My Cluster Cmd 35 gzipDF = sql("SELECT * FROM fireCallsCSVgzip") 3 %timeit -n1 -r1 gzipDF.select("City").where("City == 'San Francisco'").count() ▶ (2) Spark Jobs ▶ ■ gzipDF: pyspark.sql.dataframe.DataFrame = [Call Number: integer, Unit ID: string ... 32 more fields] 24.7 s \pm 0 ns per loop (mean \pm std. dev. of 1 run, 1 loop each) Command took 24.95 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 7:52:07 PM on My Cluster ► I bzipDF: pyspark.sql.dataframe.DataFrame = [Call Number: integer, Unit ID: string ... 32 more fields] lmin 19s \pm 0 ns per loop (mean \pm std. dev. of 1 run, 1 loop each) Command took 1.33 minutes -- by SMCPRO RTY9GMAIL.COM at 6/29/2021, 7:52:08 PM on My Cluster Cmd 37 Reading from Parquet Files "Apache Parquet is a columnar storage format available to any project in the Hadoop ecosystem, regardless of the choice of data processing framework, data model or programming language." Parquet https://parquet.apache.org **About Parquet Files** Free & Open Source Increased query performance over row-based data stores. · Provides efficient data compression. Designed for performance on large data sets. · Supports limited schema evolution Is a splittable "file format" · A Column-Oriented data store

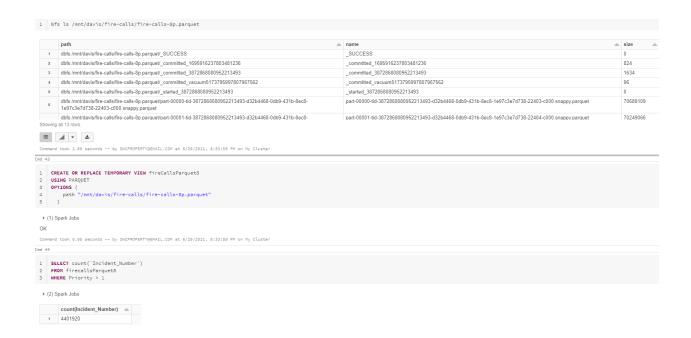
Row Format			Column Format						
ID	Name	Score	ID:	1	2	3			
1	john	4.1	Name:	john	mike	sally			
2	mike	3.5	Score:	4.1	3.5	6.4			
3	sally	6.4							

Suhaimi William Chan June 27, 2021 Page | 18

See also

- https://parquet.apache.org
- · https://en.wikipedia.org/wiki/Apache_Parquet





Command took 3.37 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 8:34:02 PM on My Cluster

Cmd 45

This file has 8 partitions rather than 1. Look at the speed improvement!

Cmd 46

Reading CSV

- spark.read.csv(..)
- . There are a large number of options when reading CSV files including headers, column separator, escaping, etc.
- . We can allow Spark to infer the schema at the cost of first reading in the entire file
- · Large CSV files should always have a schema pre-defined

Cmd 47

Reading Parquet

- spark.read.parquet(..)
- · Parquet files are the preferred file format for big-data
- . It is a columnar file format
- . It is a splittable file format
- . It offers a lot of performance benefits over other formats including predicate push down
- . Unlike CSV, the schema is read in, not inferred
- · Reading the schema from Parquet's metadata can be extremely efficient

Cmd 48

Comparison

Туре	Inference Type	Inference Speed	Reason	Should Supply Schema?
CSV	Full-Data-Read	Slow	File size	Yes
Parquet	Metadata-Read	Fast/Medium	Number of Partitions	No (most cases)
Tables	n/a	n/a	Predefined	n/a

JSON	Full-Read-Data	Slow	File size	Yes
Text	Dictated	Zero	Only 1 Column	Never
JDBC	DB-Read	Fast	DB Schema	No

Cmd 49

Reading Tables

- spark.read.table(..)
- . The Databricks platform allows us to register a huge variety of data sources as tables via the Databricks UI
- . Any DataFrame (from CSV, Parquet, whatever) can be registered as a temporary view
- Tables/Views can be loaded via the DataFrameReader to produce a DataFrame
- · Tables/Views can be used directly in SQL statements

Cmd 50

Reading JSON

- spark.read.json(..)
- . JSON represents complex data types unlike CSV's flat format
- · Has many of the same limitations as CSV (needing to read the entire file to infer the schema)
- · Like CSV has a lot of options allowing control on date formats, escaping, single vs. multiline JSON, etc.

Cmd 51

Reading Text

- spark.read.text(..)
- Reads one line of text as a single column named value
- . Is the basis for more complex file formats such as fixed-width text files

Reading JDBC

- spark.read.jdbc(..)
- · Requires one database connection per partition
- · Has the potential to overwhelm the database
- · Requires specification of a stride to properly balance partitions

Cmd 53

© 2020 Databricks, Inc. All rights reserved.

Apache, Apache Spark, Spark and the Spark logo are trademarks of the Apache Software Foundation.

Privacy Policy | Terms of Use | Support

3.5 Schemas and Types

☆ In this notebook you:

- Motivate the use of schemas and ty
- Read from JSON without a schema
- · Read from JSON with a schema



☆ Why Schemas Matter

Schemas are at the heart of data structures in Spark. A schema describes the structure of your data by naming columns and declaring the type of data in that column. Rigorously enforcing schemas leads to significant performance optimizations and reliability of code.

Why is open source Spark so fast, and why is Databricks Runtime even faster? While there are many reasons for these performance improvements, two key reasons are:

- · First and foremost, Spark runs first in memory rather than reading and writing to disk
- Second, using DataFrames allows Spark to optimize the execution of your queries because it knows what your data looks like

Two pillars of computer science education are data structures, the organization and storage of data and algorithms, and the computational procedures on that data. A rigorous understanding of computer science involves both of these domains. When you apply the most relevant data structures, the algorithms that carry out the computation become significantly more eloquent

Schemas with Semi-Structured JSON Data

Tabular data, such as that found in CSV files or relational databases, has a formal structure where each observation, or row, of the data has a value (even if it's a NULL value) for each feature, or column, in the data set.

Semi-structured data does not need to conform to a formal data model. Instead, a given feature may appear zero, once, or many times for a given observation.

Semi-structured data storage works well with hierarchical data and with schemas that may evolve over time. One of the most common forms of semi-structured data is JSON data, which consists of attribute-value pairs.

🛱 Reading from JSON w/ InferSchema

Reading in JSON isn't that much different than reading in CSV files.

Let's start with taking a look at all the different options that go along with reading in JSON files.

JSON Lines

Much like the CSV reader, the JSON reader also assumes:

- . That there is one JSON object per line and
- That it's delineated by a new-line

mat is referred to as JSON Lines or newline-delimited JSON

More information about this format can be found at http://isonlines.org.

Note: Spark 2.2 was released on July 11th 2016. With that comes File 10 improvements for CSV & JSON, but more importantly, support for parsing multi-line JSON and CSV files. You can read more about that (and other features in Spark 2.2) in the Databricks Blog

Cnd 9

Take a look at the sample of our JSON data

1 %fs ls /mnt/davis/fire-calls/fire-calls-truncated.json



Like we did with the CSV file, we can use %fs head ... to take a look at the first few lines of the file

1 %fs head /mnt/davis/fire-calls/fire-calls-truncated.json

Truncated to first 65536 bytes]

("Call Number:1009118," Unit D':"100", "Incident Number:10625," Call Type"; "Medical Incident", "Call Date"; "04/12/2000", "Netch Date"; "04/12/2000", "Received Date"; "04/12/2000 09:27:48 PM", "Entry Date"; "04/12/2000 09:27:48 PM", "Dispatch Date"; "04/12/2000 09:27:48 PM", "Entry Date"; "04/12/2000 09:27:48 PM", "Dispatch Date"; "04/12/2000 09:29:12 PM", "Nesponse Date"; "04/12/2000 09:39:12 PM", "Nesponse Date"; "04/12/2000 09:39:13 PM", "05/12/2000 09:39:14 PM", "05/12/2000 09:39:15 PM", "05/12/200

Read the JSON File

The command to read in JSON looks very similar to that of CSV.

In addition to reading the JSON file, we will also print the resulting schema.

```
Cmd 14
     CREATE OR REPLACE TEMPORARY VIEW fireCallsJSON
 2
     USING JSON
    OPTIONS (
          path "/mnt/davis/fire-calls/fire-calls-truncated.json"
        )
 5
   ▶ (1) Spark Jobs
 ΟK
 Command took 9.92 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 9:07:39 PM on My Cluster
Cmd 15
 1
      DESCRIBE fireCallsJSON
          col_name
                                                 data_type
                                                                 comment
          ALS Unit
                                                 boolean
                                                                 null
          Address
                                                 string
                                                                 null
         Available DtTm
                                                 string
                                                                 null
         Battalion
                                                 string
                                                                 null
          Box
                                                 string
                                                                 null
          Call Date
                                                 string
                                                                 null
```

Showing all 34 rows.



7 Call Final Disposition

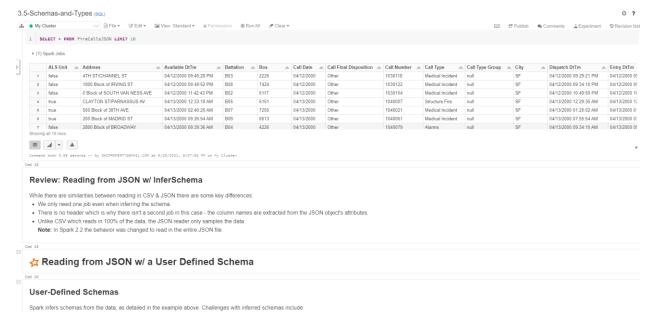
Command took 0.09 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 9:07:53 PM on My Cluster

strina

null

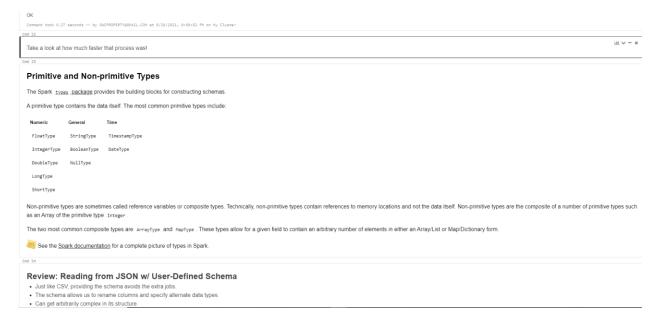
Cmd 16

Take a look at the table.



- · Schema inference means Spark scans all of your data, creating an extra job, which can affect performance
- Consider providing alternative data types (for example, change a Long to a Integer)
- . Consider throwing out certain fields in the data, to read only the data of interest

```
CREATE OR REPLACE TEMPORARY VIEW fireCallsJSON (
2
      `Call Number` INT,
      'Unit ID' STRING,
3
4
      'Incident Number' INT,
      'Call Type' STRING,
5
      'Call Date' STRING,
6
7
      'Watch Date' STRING,
      'Received DtTm' STRING,
8
      `Entry DtTm` STRING,
9
      'Dispatch DtTm' STRING,
10
      `Response DtTm` STRING,
11
12
      'On Scene DtTm' STRING,
      `Transport DtTm` STRING,
13
      `Hospital DtTm` STRING,
14
      'Call Final Disposition' STRING,
15
      `Available DtTm` STRING,
16
17
      `Address` STRING,
      'City' STRING,
18
      'Zipcode of Incident' INT,
19
      `Battalion` STRING,
20
      'Station Area' STRING,
21
22
      'Box' STRING,
      'Original Priority' STRING,
23
      `Priority` STRING,
24
      `Final Priority` INT,
25
      'ALS Unit' BOOLEAN,
26
27
      'Call Type Group' STRING,
      `Number of Alarms` INT,
28
      'Unit Type' STRING,
29
      'Unit sequence in call dispatch' INT,
30
      `Fire Prevention District` STRING,
31
32
      `Supervisor District` STRING,
      'Neighborhooods - Analysis Boundaries' STRING,
33
      `Location` STRING,
34
      'RowID' STRING
35
36
37
    USING JSON
38
    OPTIONS (
        path "/mnt/davis/fire-calls/fire-calls-truncated.json"
39
40
```



3.6 Writing Data

☆ In this notebook you:

Employ the writing design pattern

Control partitions for database writes using REPARTITION and COALESCE hints

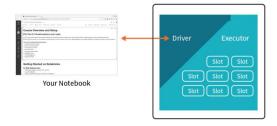


ind 5

Writing to a database in Spark differs from other tools largely due to its distributed nature. There are a number of variables that can be tweaked to optimize performance, largely relating to how data is organized on the cluster. Partitions are the first step in understanding performant database connections.

A partition is a portion of your total data set, which is divided into many of these portions so Spark can distribute your work across a cluster.

The other concept needed to understand Spark's computation is a slot (also known as a core). A slot/core is a resource available for the execution of computation in parallel. In brief, a partition refers to the distribution of data while a slot refers to the distribution of computation.



As a general rule of thumb, the number of partitions should be a multiple of the number of cores. For instance, with 5 partitions and 8 slots, 3 of the slots will be underutilized. With 9 partitions and 8 slots, a job will take twice as long as it waits for the extra partition to finish.

```
Import the dataset.

Tod 7

CREATE OR REPLACE TEMPORARY VIEW fireCallsCSV

SING CSV

OPTIONS (

path */mot/davis/fire-calls-truncated.csv",

header "true",

sep ":"

) (1) Spark Jobs

OK
```

Command took 0.84 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 9:25:22 PM on My Cluster Switch to the Python API to run write commands. You can do this %python and then run sql(" <My Normal SQL Query> ") 1 %python 2 df = sql("SELECT * FROM fireCallsCSV") 3 4 display(df) ▶ (1) Spark Jobs → ■ df: pyspark.sql.dataframe.DataFrame = [Call Number: string, Unit ID: string ... 32 more fields]
 Received DrTm
 Entry DrTm
 Dispatch DrTm
 Response DrTm
 On Scene DrTm

 04/12/2000 09.27.45 PM
 04/12/2000 09.28.58 PM
 04/12/2000 09.29.21 PM
 04/12/2000 09.31.26 PM
 04/12/2000 09.33.49 PM

 04/12/2000 09.31.55 PM
 04/12/2000 09.33.48 PM
 04/12/2000 09.34.10 PM
 04/12/2000 09.35.59 PM
 04/12/2000 09.35.59 PM
 04/12/2000 09.35.59 PM
 # L | + | + | nds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 9:27:02 PM on My Cluster Cmd 10 Write using the .write(" < path> ") method on the DataFrame Cmd 11 1 %python 2 df.write.mode("OVERWRITE").csv(username + "/fire-calls.csv") Command took 16.32 seconds -- by SWCPROPERTYAGNAIL.COM at 6/29/2021, 9:27:05 PM on My Clust-Take a look at the file we just wrote. 08-1-c000.csv', size-11655265),
FileInfo(path='dbfs:/SMCPROPERTY@GMAIL.COM/fire-calls.csv/part-00002-tid-469682914843696193-c27ede8b-f5a1-4feb-87bb-fcb846eb7ac8-1 FileIn(gath-'dbfa:/MCPOERTY@MIL.COM/fire-calls.csv/part-00005-tid-469682914843096139-c27edebb-fs1-4feb-77bb-fcb046eb7ac3-112-1-000.csv', name-'part-00005-tid-469682914843096133-c27edebb-fs1-4feb-77bb-fcb046eb7ac3-113-1-000.csv', name-'part-00006-tid-469682914843096133-c27edebb-fs1-4feb-87bb-fcb046eb7ac3-113-1-000.csv', name-'part-00006-tid-469682914843096133-c27edebb-fs1-4feb-87bb-fcb046eb7ac3-113-1-000.csv', name-'part-00006-tid-469682914843096133-c27edebb-fs1-4feb-87bb-fcb046eb7ac3-113-1-000.csv', name-'part-00006-tid-469682914843096133-c27edebb-fs2b-14feb-87bb-fcb046eb7ac3-114-1-000.csv', name-'part-00007-tid-469682914843096133-c27edebb-fs2b-4feb-87bb-fcb046eb7ac3-114-1-000.csv', name-'part-00007-tid-46968291484309133-c27edebb-fs2b Notice that it wrote a number of different parts, one for each partition. Now let's get a sense for how to check how many partitions our data has, which once again is in Python. 1 %python 2 df.rdd.getNumPartitions() Out[5]: 8

☆ Controlling Concurrency

Cmd 17

In the context of JDBC database writes, the number of partitions determine the number of connections used to push data through the JDBC API. There are two ways to control this parallelism:

Hint	Transformation Type	Use	Evenly distributes data across partitions?
SELECT /*+ COALESCE(n) */	narrow (does not shuffle data)	reduce the number of partitions	no
SELECT /*+ REPARTITION(n) */	wide (includes a shuffle operation)	increase the number of partitions	yes

Cmd 18

Create a new temporary view that coalesces all of our data to a single partition. This can be done using a hint.

Now check the number of partitions.

We now have a single partition. What if we want more partitions? Let's use a REPARTITION hint to evenly distribute our data across 8 partitions.

```
CREATE OR REPLACE TEMPORARY VIEW fireCallsCSV8p

AS
SELECT /*+ REPARTITION(8) */ *
FROM fireCallsCSV
```

OK

Command took 0.18 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 9:42:57 PM on My Cluster

1 | %python 2 | sql("SELECT * FROM fireCallsCSV8p").rdd.getNumPartitions()

• (1) Spark Jobs
Out[7]: 8
Command took 6.82 seconds -- by SMCPROPERTYGGMAIL.COM at 6/29/2021, 9:42:58 PM on My Cluster

nd 25

Now you can save the results.

```
1 %python
2 sql("SELECT * FROM fireCallsCSV8p").write.mode("OVERWRITE").csv(username + "/fire-calls-repartitioned.csv")

• (2) Spark Jobs

Command took 14.90 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 9:43:35 PM on My Cluster
```

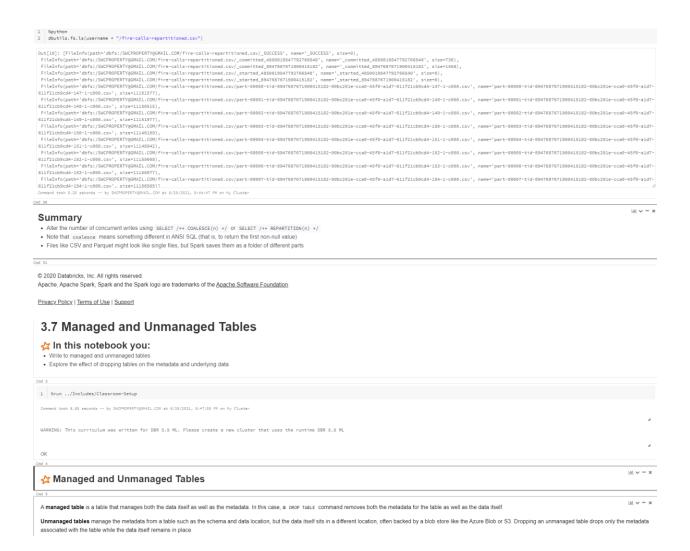
Cmd 27

Click the arrow next to Spark Jobs under the following code cell in order to see a breakdown of the job you triggered. Click the next arrow to see a breakdown of the stages.

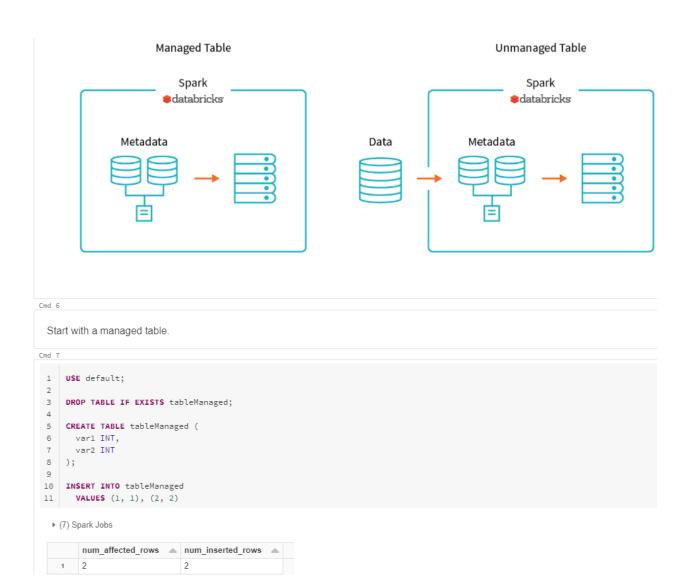
When you repartitioned the DataFrame to 8 partitions, 8 stages were needed, one to write each partition of the data.

Cmd 28

View the results.

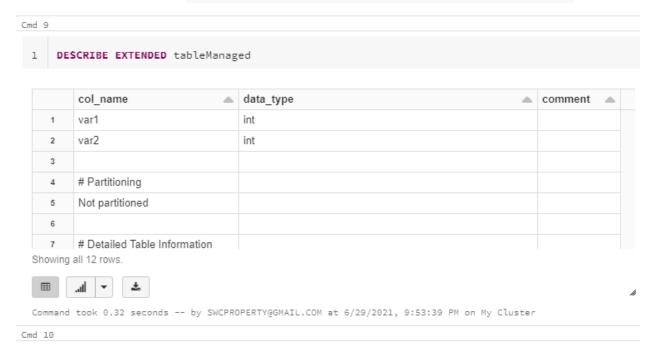


June 27, 2021 Suhaimi William Chan Page | 30



Use DESCRIBE EXTENDED to describe the contents of the table. Scroll down to see the table Type.

Notice the location is also dbfs:/user/hive/warehouse/< your database >/tablemanaged .



Now use an external, or unmanaged, table

```
DROP TABLE IF EXISTS tableUnmanaged;

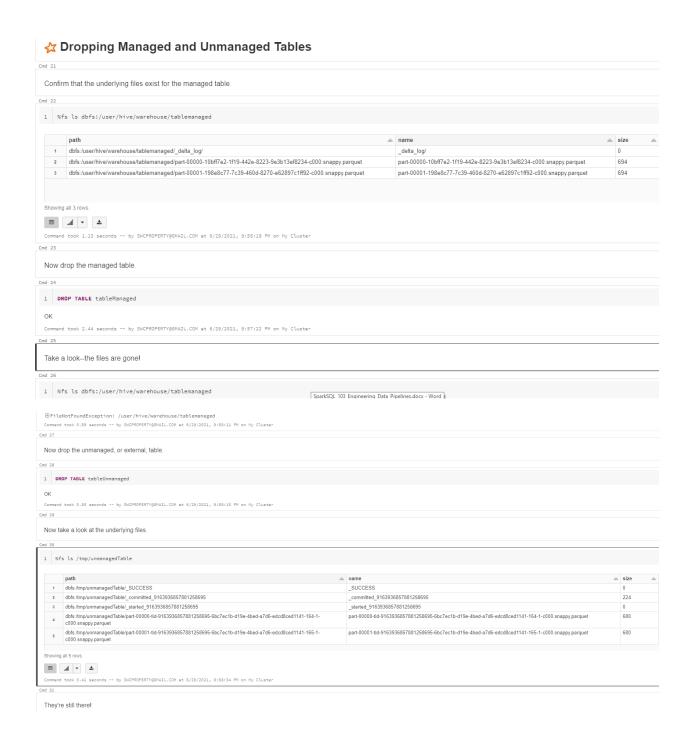
CREATE EXTERNAL TABLE tableUnmanaged (
varl INT,
var2 INT

)

STORED AS parquet
LOCATION '/tmp/unmanagedTable'
```

Command took 0.37 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 9:53:44 PM on My Cluster

Describe the table and look for the Type Cmd 13 1 DESCRIBE EXTENDED tableUnmanaged col_name data_type comment var1 int null 2 var2 int null 4 # Detailed Table Information 5 Database default Table tableunmanaged 7 Owner Showing all 19 rows. .il - ± \blacksquare Command took 0.11 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 9:53:48 PM on My Cluster Cmd 14 This is an external, or managed table. If we were to shut down our cluster, this data will persist. Now insert values into the table. Cmd 15 1 INSERT INTO tableUnmanaged 2 VALUES (1, 1), (2, 2) ▶ (1) Spark Jobs OK Command took 2.53 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 9:54:59 PM on My Cluster Take a look at the result. Cmd 17 1 SELECT * FROM tableUnmanaged ▶ (2) Spark Jobs Showing all 2 rows. ± ... = Now view the underlying files in where the data was persisted. 1 %fs ls /tmp/unmanagedTable path
dbfs:/tmp/unmanagedTable/_SUCCESS _SUCCESS 2 dbfs:/tmp/unmanagedTable/_committed_9163936857881258695 3 dbfs:/tmp/unmanagedTable/_started_9163936857881258695 committed 9163936857881258695 224 _started_9163936857881258695 dbfs/tmplummanagedTable/part-00000-tid-9163936857881258695-6bc7ec1b-d19e-4bed-a7d6-edcd8ced1141-164-1-c000 anapyr parquet part-00000-tid-9163936857881258695-6bc7ec1b-d19e-4bed-a7d6-edcd8ced1141-164-1-c000.snappy.parquet $\label{eq:dbs} $$dbs:/tmp/unmanagedTable/part-00001-tid-9163936857881258695-6bc7ec1b-d19e-4bed-a7d6-edcd8ced1141-165-1c000.snappy.parquet$ part-00001-tid-9163936857881258695-6bc7ec1b-d19e-4bed-a7d6-edcd8ced1141-165-1-c000.snappy.parquet 680 Showing all 5 rows. ± + h. = Command took 1.37 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 9:56:12 PM on My Cluster



June 27, 2021 Suhaimi William Chan Page | 34

Summary

- Use external/unmanaged tables when you want to persist your data once the cluster has shut down
- · Use managed tables when you only want ephemeral data

Cmd 33

© 2020 Databricks, Inc. All rights reserved.

Apache, Apache Spark, Spark and the Spark logo are trademarks of the Apache Software Foundation.

Privacy Policy | Terms of Use | Support

Engineering Data Pipelines

Module 3 Assignment

☆ In this assignment you:

- Create a table with persistent data and a specified schema
- · Populate table with specific entries
- Change partition number to compare query speeds



Check that the data type of each column is what we want.

Question 1

What type of table is newTable ? "EXTERNAL" or "MANAGED"?

```
1 DESCRIBE EXTENDED newTable
```

	col_name	data_type	comment	
12	Created By	Spark 3.1.0		
13	Туре	EXTERNAL		
14	Provider	hive		
15	Table Properties	[transient_lastDdlTime=1625030369]		
16	Location	dbfs:/tmp/newTableLoc		
17	Serde Library	org.apache.hadoop.hive.ql.io.parquet.serde.ParquetHiveSerDe		

Run the following cell to read in the data stored at /mnt/davis/fire-calls/fire-calls-truncated.json . Check that the columns of the data are of the correct types (not all strings).

```
1 CREATE OR REPLACE TEMPORARY VIEW fireCallsJSON (
            `ALS Unit' boolean,
'Address' string,
'Available DtTm' string,
'Battalion' string,
             'Box' string,
'Call Date' string,
'Call Final Disposition' string,
             'Call Number' long,
'Call Type' string,
             'Call Type Group' string,
'City' string,
'Dispatch DtTm' string,
 12
            `Entry DtTm` string,
`Final Priority` long,
            Fine Prevention District' string,

'Hospital DtTm' string,

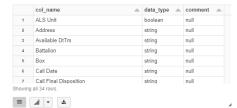
'Incident Number' long,
'Location' string,

'Neighborhooods - Analysis Boundaries' string,

'Number of Alarms' long,
'O. Scene DtTm' string,
19
21
             'On Scene DtTm' string,
            Original Priority' string,
'Priority' string,
'Received DtTm' string,
'Response DtTm' string,
23
24
25
26
             `RowID` string,
            'Station Area' string,
'Supervisor District' string,
28
             `Transport DtTm` string,
`Unit ID` string,
 30
31
32
            'Unit Type' string,
'Unit sequence in call dispatch' long,
            `Watch Date` string,

`Zipcode of Incident` long
35
36 )
37 USING JSON
38 OPTIONS (
              path "/mnt/davis/fire-calls/fire-calls-truncated.json"
39
40 );
```

42 DESCRIBE fireCallsJSON



Command took 0.21 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 10:23:09 PM on My Cluster

Take a look at the table to make sure it looks correct.

Cmd 12

1 SELECT * FROM fireCallsJSON LIMIT 10

▶ (1) Spark Jobs

		ALS Unit 🔺	Address	Available DtTm	Battalion 🔺	Box 📤	Call Date 🔺	Call Final Disposition 🔺	Call Number 🔺	Call Type	Call Type Group 🔺	City
	1	false	4TH ST/CHANNEL ST	04/12/2000 09:45:28 PM	B03	2226	04/12/2000	Other	1030118	Medical Incident	null	SF
	2	false	1800 Block of IRVING ST	04/12/2000 09:49:52 PM	B08	7424	04/12/2000	Other	1030122	Medical Incident	null	SF
	3	false	0 Block of SOUTH VAN NESS AVE	04/12/2000 11:42:43 PM	B02	5117	04/12/2000	Other	1030154	Medical Incident	null	SF
	4	true	CLAYTON ST/PARNASSUS AV	04/13/2000 12:33:18 AM	B05	5151	04/13/2000	Other	1040007	Structure Fire	null	SF
	5	true	500 Block of 38TH AVE	04/13/2000 02:40:25 AM	B07	7255	04/13/2000	Other	1040021	Medical Incident	null	SF
	6	true	200 Block of MADRID ST	04/13/2000 09:26:54 AM	B09	0613	04/13/2000	Other	1040061	Medical Incident	null	SF
	7	false	2800 Block of BROADWAY	04/13/2000 09:39:36 AM	B04	4226	04/13/2000	Other	1040079	Alarms	null	SF
S	howing	all 10 rows.										

Lal + ±

Command took 0.50 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 10:23:16 PM on My Cluster Command took 0.50 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 10:23:16 PM on My Cluster Command took 0.50 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 10:23:16 PM on My Cluster Command took 0.50 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 10:23:16 PM on My Cluster Command took 0.50 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 10:23:16 PM on My Cluster Command took 0.50 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 10:23:16 PM on My Cluster Command took 0.50 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 10:23:16 PM on My Cluster Command took 0.50 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 10:23:16 PM on My Cluster Command took 0.50 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 10:23:16 PM on My Cluster Command took 0.50 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 10:23:16 PM on My Cluster Command to 6/29/2021, 10:23:16/2021, 10:23:16/2021

Now let's populate newTable with some of the rows from the fireCallsJSON table you just loaded. We only want to include fire calls whose Final Priority is 3.



▶ (2) Spark Jobs

	Address	City	Battalion 🔺	Box 🔺
1	2500 Block of BROADWAY	SF	B04	4122
2	100 Block of BERRY ST	SF	B03	2171
3	100 Block of APPLETON AVE	SF	B06	5646
4	400 Block of JONES ST	SF	B03	1461
5	300 Block of GRANT AVE	SF	B01	1315
6	6TH ST/MARKET ST	SF	B03	2248
7	800 Block of 3RD ST	SF	B03	2150

Truncated results, showing first 1000 rows.



Command took 14.28 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 10:31:57 PM on My Cluster

Question 2

How many rows are in newTable ?

1 select count(*) from newTable

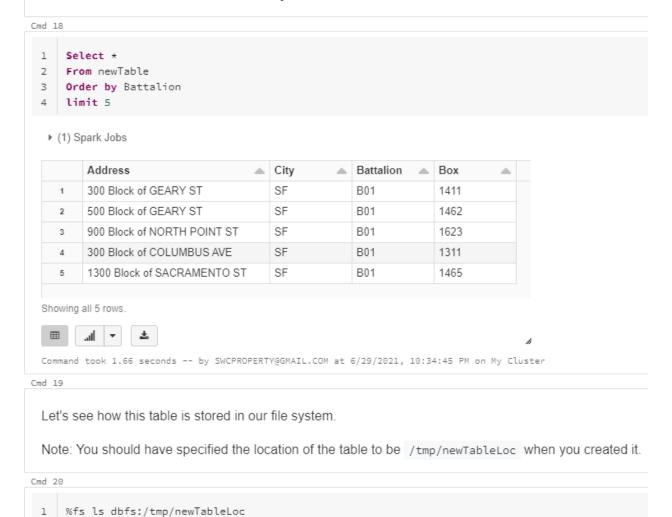
▶ (2) Spark Jobs

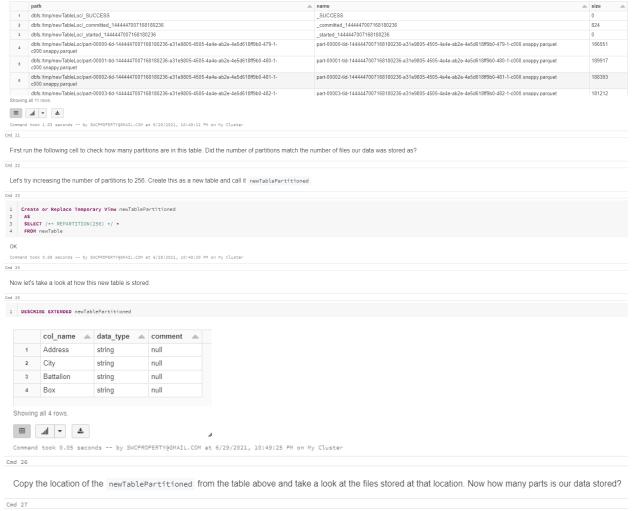


Sort the rows of newTable by ascending Battalion .

Question 3

What is the "Battalion" of the first entry in the sorted table?





27

1 %fs ls dbfs:/user/hive/warehouse/databricks.db/newtablepartitioned

FileNotFoundException: /user/hive/warehouse/databricks.db/newtablepartitioned

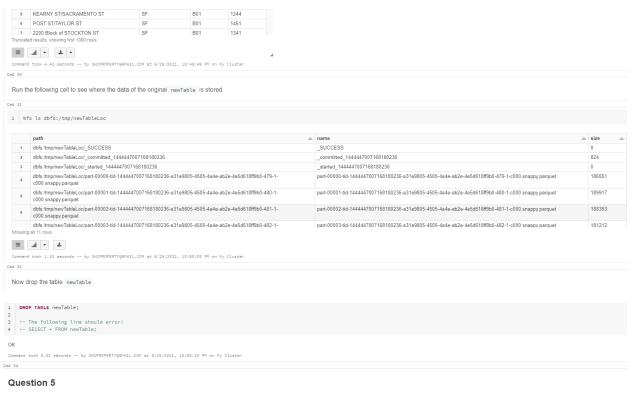
Command took 0.46 seconds -- by SWCPROPERTY@GMAIL.COM at 6/29/2021, 10:49:36 PM on My Cluster

Now sort the rows of newTablePartitioned by ascending Battalion and compare how long this query takes.

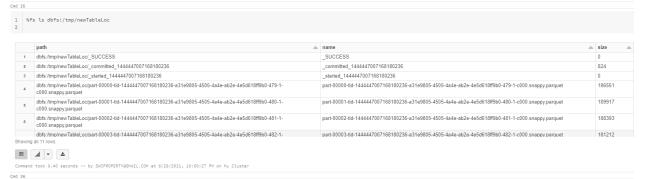
Question 4

Was this query faster or slower on the table with increased partitions?





Does the data stored within the table still exist at the original location (dbfs:/tmp/newTableLoc) after you dropped the table? (Answer "yes" or "no")



© 2020 Databricks, Inc. All rights reserved.

Apache, Apache Spark, Spark and the Spark logo are trademarks of the <u>Apache Software Foundation</u>