



ScholarNest

# Spark Azure Databricks

Databricks Spark Certification and beyond

Instructor: Prashant Kumar Pandey



Absolute Beginner to Specialization in Apache Spark and Azure Databricks





# Creating Spark Session

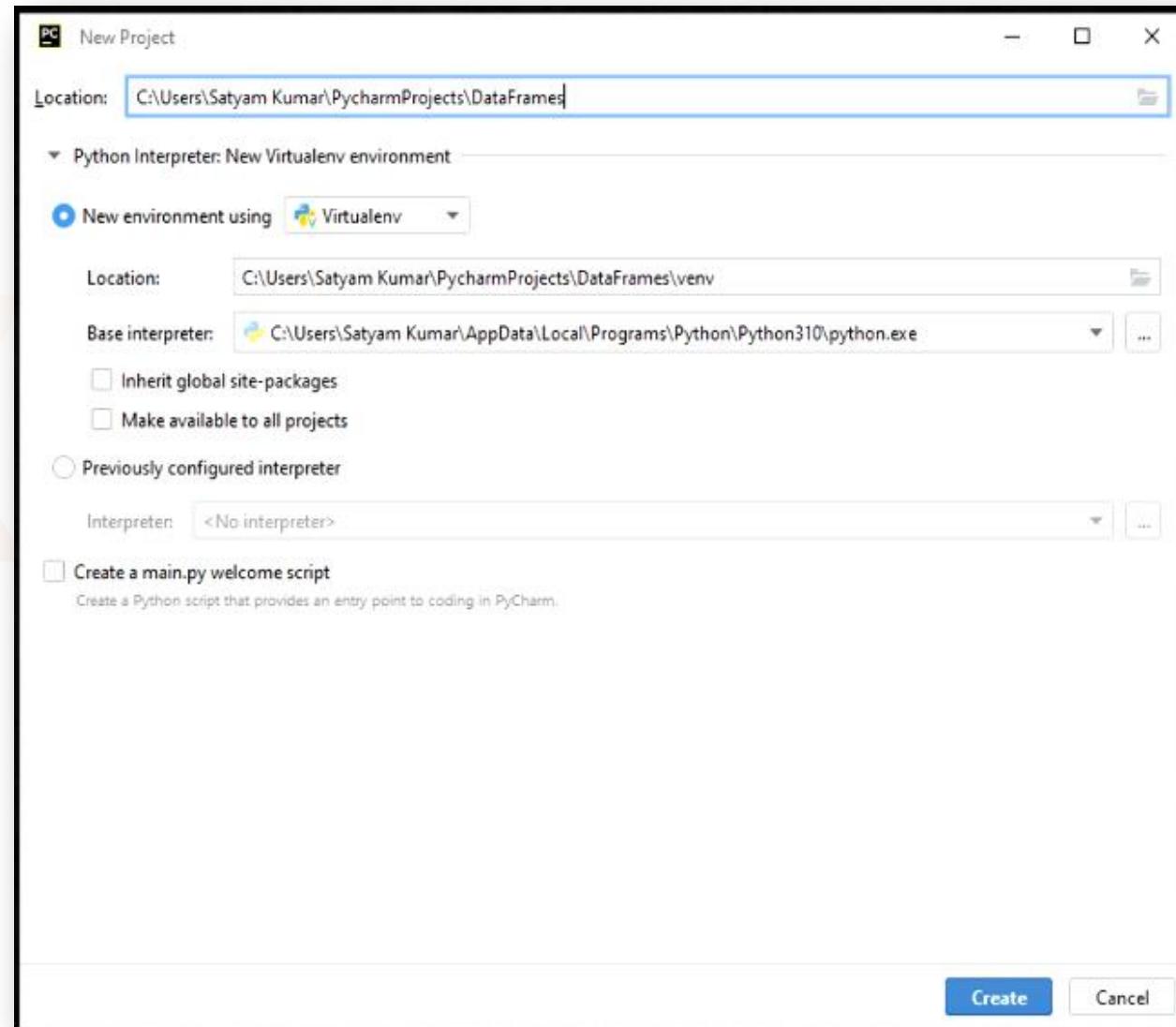
Assume you are given a task to create a small Spark job to do the following:

## Requirement

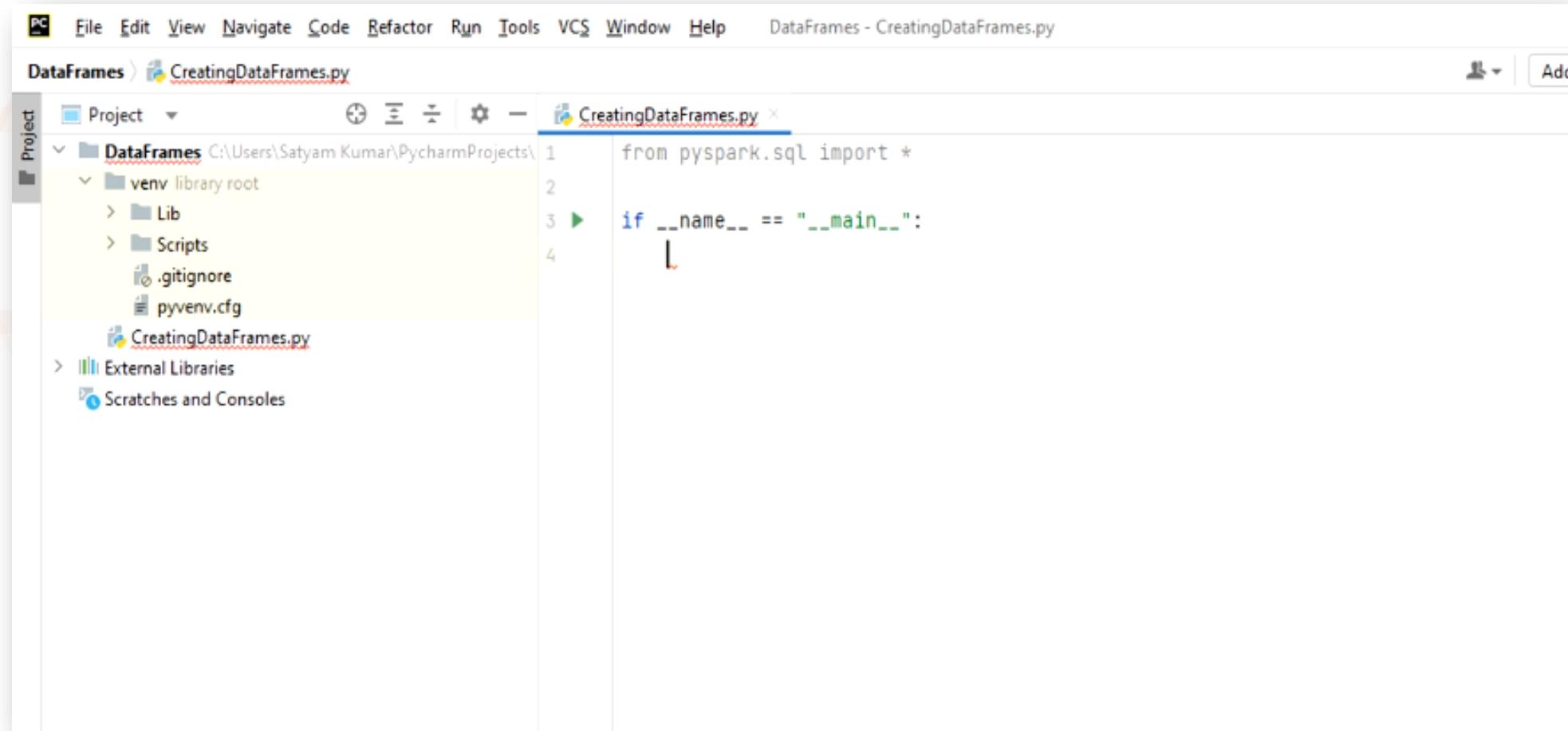
---

1. Read a CSV data file
2. Perform some standard data validation and data type correction to prepare it for analysis
3. Save the valid data using the parquet data file format

Start your PyCharm IDE and create a new project. Give a name to your project and create it. **(Reference : DataFrames)**



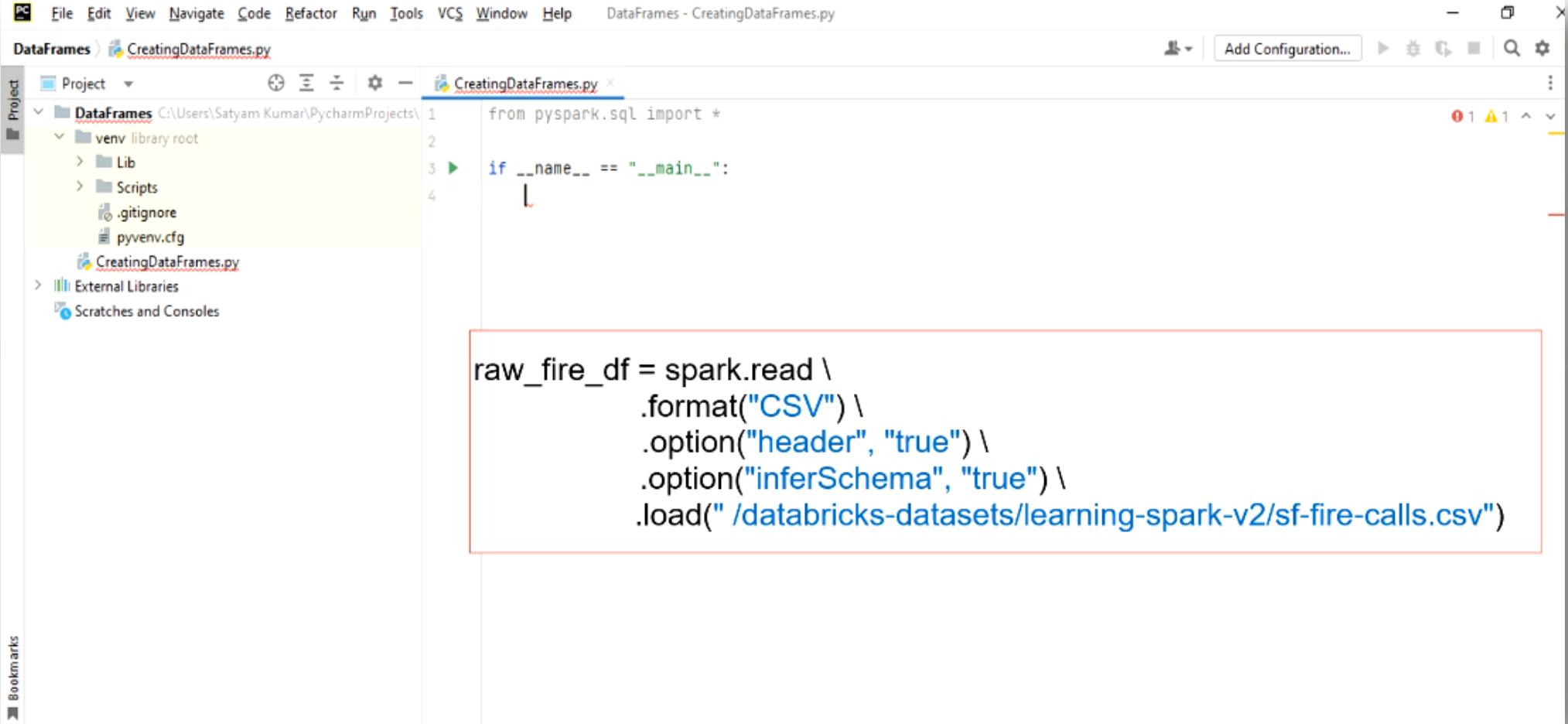
Create a new Python file in your project. We will be using the PySpark SQL package so let me import the package. I am starting a new Spark project, so I need an entry point for my application. So I have defined the main method as shown below.



The screenshot shows the PyCharm IDE interface. The title bar reads "DataFrames - CreatingDataFrames.py". The menu bar includes File, Edit, View, Navigate, Code, Refactor, Run, Tools, VCS, Window, Help. The toolbar has icons for New, Open, Save, Cut, Copy, Paste, Find, Replace, and Run. The Project tool window on the left shows a "DataFrames" project structure with a "venv" folder containing "Lib", "Scripts", ".gitignore", and "pyvenv.cfg". A file named "CreatingDataFrames.py" is selected. The code editor window contains the following Python code:

```
from pyspark.sql import *
if __name__ == "__main__":
    
```

The first requirement is to read a CSV data file. How will you do it? Here it is. We used this code to read a CSV file and create a Spark DataFrame. So we can use the same code template here.



The screenshot shows the PyCharm IDE interface. The top menu bar includes File, Edit, View, Navigate, Code, Refactor, Run, Tools, VCS, Window, and Help. The title bar says "DataFrames - CreatingDataFrames.py". The left sidebar shows a project structure with a "DataFrames" folder containing "CreatingDataFrames.py", "venv", "Lib", "Scripts", ".gitignore", and "pyenv.cfg". The main editor window displays the following Python code:

```
from pyspark.sql import *
if __name__ == "__main__":
    raw_fire_df = spark.read \
        .format("CSV") \
        .option("header", "true") \
        .option("inferSchema", "true") \
        .load("/databricks-datasets/learning-spark-v2/sf-fire-calls.csv")
```

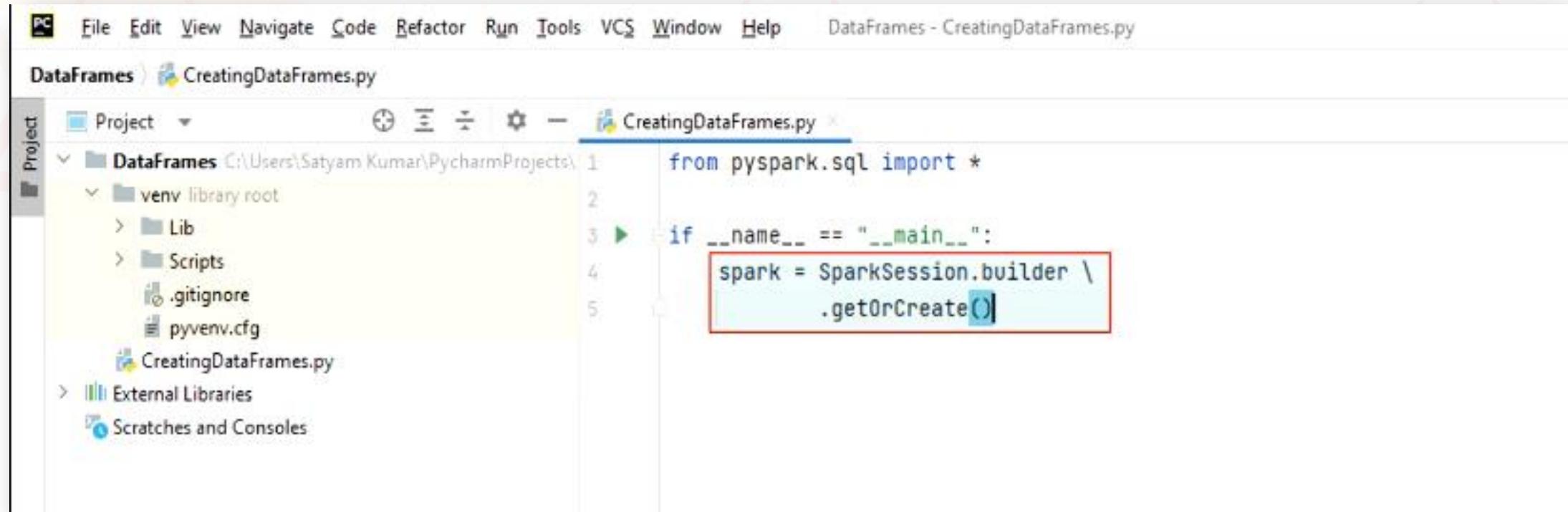
But we have a small problem. We do not have the Spark Session variable. I created the earlier example in a Databricks Notebook. The notebook environment automatically creates a Spark session object and makes it available for us. But when you are creating a Spark program in a Python file, you must create the SparkSession yourself as shown below.

The screenshot shows the PyCharm IDE interface. The top menu bar includes File, Edit, View, Navigate, Code, Refactor, Run, Tools, VCS, Window, Help, and the current file name 'DataFrames - CreatingDataFrames.py'. Below the menu is a toolbar with icons for Project, New, Open, Save, and Run. The left sidebar is titled 'Project' and shows a tree structure for the 'DataFrames' project. It includes a 'venv' folder containing 'Lib', 'Scripts', '.gitignore', and 'pyenv.cfg', and a file named 'CreatingDataFrames.py'. The main editor window displays the following Python code:

```
from pyspark.sql import *
if __name__ == "__main__":
    spark = SparkSession.builder \
        .getOrCreate()
```

The line 'spark = SparkSession.builder \ .getOrCreate()' is highlighted with a red rectangular box.

But we have a small problem. We do not have the Spark Session variable. I created the earlier example in a Databricks Notebook. The notebook environment automatically creates a Spark session object and makes it available for us. But when you are creating a Spark program in a Python file, you must create the SparkSession yourself as shown below. The code shown here is the bare minimum to start writing the Spark program in a python file.



The screenshot shows the PyCharm IDE interface. The top menu bar includes File, Edit, View, Navigate, Code, Refactor, Run, Tools, VCS, Window, and Help. The current file is "DataFrames - CreatingDataFrames.py". The Project tool window on the left shows a "DataFrames" project with a "venv" library root containing "Lib", "Scripts", ".gitignore", and "pyenv.cfg". The "CreatingDataFrames.py" file is open in the editor. The code is as follows:

```
from pyspark.sql import *
if __name__ == "__main__":
    spark = SparkSession.builder \
        .getOrCreate()
```

The line "spark = SparkSession.builder \ .getOrCreate()" is highlighted with a red rectangular box.

Spark Session is your entry point for Spark programming, and it offers you the following attributes and methods.

## SparkSession

### Attributes

<code>builder</code>	A class attribute having a <code>Builder</code> to construct <code>SparkSession</code> instances.
<code>read</code>	Returns a <code>DataFrameReader</code> that can be used to read data in as a <code>DataFrame</code> .
<code>readStream</code>	Returns a <code>DataStreamReader</code> that can be used to read data streams as a streaming <code>DataFrame</code> .
<code>catalog</code>	Interface through which the user may create, drop, alter or query underlying databases, tables, functions, etc.
<code>conf</code>	Runtime configuration interface for Spark
<code>udf</code>	Returns a <code>UDFRegistration</code> for UDF registration
<code>sparkContext</code>	Returns the underlying <code>SparkContext</code> .
<code>version</code>	The version of Spark on which this application is running.

### Methods

<code>createDataFrame(data[, schema, ...])</code>	Creates a <code>DataFrame</code> from an <code>RDD</code> , a list or a <code>pandas.DataFrame</code> .
<code>getActiveSession()</code>	Returns the active <code>SparkSession</code> for the current thread, returned by the builder
<code>newSession()</code>	Returns a new <code>SparkSession</code> as new session, that has separate <code>SQLConf</code> , registered temporary views and UDFs, but shared <code>SparkContext</code> and table cache.
<code>range(start[, end, step, numPartitions])</code>	Create a <code>DataFrame</code> with single <code>pyspark.sql.types.LongType</code> column named <code>id</code> , containing elements in a range from <code>start</code> to <code>end</code> (exclusive) with step value <code>step</code> .
<code>sql(sqlQuery)</code>	Returns a <code>DataFrame</code> representing the result of the given query.
<code>stop()</code>	Stop the underlying <code>SparkContext</code> .
<code>table(tableName)</code>	Returns the specified table as a <code>DataFrame</code> .

Let's focus on the attributes for now. The first one is the builder. Once you have the session builder, you can use builder methods to create and configure your spark session. And here is the list of session builder methods. So we have five session builder methods. The appName allows you to set a unique name for your Spark application.

SparkSession	
<strong>Attributes</strong>	
<code>builder</code>	A class attribute having a <code>Builder</code> to construct <code>SparkSession</code> instances
<code>read</code>	Returns a <code>DataFrameReader</code> that can be used to read data in as a <code>DataFrame</code> .
<code>readStream</code>	Returns a <code>DataStreamReader</code> that can be used to read data streams as a streaming <code>DataFrame</code> .
<code>catalog</code>	Interface through which the user may create, drop, alter or query underlying databases, tables, functions, etc.
<code>conf</code>	Runtime configuration interface for Spark.
<code>udf</code>	Returns a <code>UDFRegistration</code> for UDF registration.
<code>sparkContext</code>	Returns the underlying <code>SparkContext</code> .
<code>version</code>	The version of Spark on which this application is running.
<strong>Builder Methods</strong>	
<code>SparkSession.builder.appName(name)</code>	Sets a name for the application, which will be shown in the Spark web UI.
<code>SparkSession.builder.master(master)</code>	Sets the Spark master URL to connect to, such as "local" to run locally, "local[4]" to run locally with 4 cores, or "spark://master:7077" to run on a Spark standalone cluster.
<code>SparkSession.builder.config([key, value, conf])</code>	Sets a config option.
<code>SparkSession.builder.enableHiveSupport()</code>	Enables Hive support, including connectivity to a persistent Hive metastore, support for Hive SerDes, and Hive user-defined functions.
<code>SparkSession.builder.getOrCreate()</code>	Gets an existing <code>SparkSession</code> or, if there is no existing one, creates a new one based on the options set in this builder.

The next one is to set the location of the master node of your Spark cluster. I do not recommend using this method at all. You should never set the master from your code.

Why? Because you will run your code on your local machine or in the development cluster. Then you want to run the same application in a testing cluster. Finally, the code goes to the production cluster.

The master node location is different everywhere. So we do not hardcode it in the application. Instead, we set the master node location from outside the application code.

The config method allows you to configure your spark session.

The next one is to enable Hive support in your Spark application. Enabling Hive support allows you to access Hive database tables from Spark. You should enable this option when your target is the Hadoop cluster. But this option is not meaningful in the Databricks Cloud environment. Because you do not have Hadoop and Hive in Databricks Cloud.

The last and the most important one is to getOrCreate a Spark session. Spark application must have only one active Spark session at any given point. So the getOrCreate API protects you from incorrectly creating two sessions. If you already created a SparkSession earlier, the getOrCreate will return the same session; else, it will create a new session. So even if you call the getOrCreate method more than once, you will get the same Spark session. Make sense?

We have quite a few other attributes of Spark Session.

The read and the readStream are the most important. They give you a DataFrameReader and open your gateway to the DataFrame API.

The catalog is your Spark metadata interface.

The conf is your spark runtime configuration interface.

The UDF here gives you a UDFRegistration interface.

The last one is the SparkContext.

**A SparkContext represents the connection to a Spark cluster and can use to create RDD and broadcast variables on the cluster.**

Memorize this definition.

The SparkContext offers you two functionalities.

1. It holds and manages a connection to your Spark cluster
2. It is your gateway to Spark Core APIs

So Spark session also offers you two functionalities.

1. It is your gateway to Spark DataFrame APIs
2. It holds the SparkContext.

Now you can combine the SparkContext into SparkSession.

1. It is your gateway to Spark DataFrame APIs
2. It holds the SparkContext.
  1. It holds and manages a connection to your Spark cluster
  2. It is your gateway to Spark Core APIs

You learned to create a SparkSession.

It is the bare minimum code. But I recommend setting a unique application name.

Simply add a line just before the getOrCreate and use Session Builder's appName() method to set a Unique application name.

The screenshot shows the PyCharm IDE interface. The top menu bar includes File, Edit, View, Navigate, Code, Refactor, Run, Tools, VCS, Window, Help, and the current file name, DataFrames - CreatingDataFrames.py. Below the menu is a toolbar with icons for Project, New, Delete, and Settings. The main window has two panes: the left pane is the Project view showing a 'DataFrames' project with a 'venv' library root containing 'Lib', 'Scripts', '.gitignore', and 'pyvenv.cfg', and a file 'CreatingDataFrames.py'; the right pane is the Code Editor with the following Python code:

```
from pyspark.sql import *
if __name__ == "__main__":
    spark = SparkSession.builder \
        .appName("CreatingDataFramesDemo") \
        .getOrCreate()
```

A blue arrow points to the line '.appName("CreatingDataFramesDemo")' in the code editor.

# Summary



You always need a SparkSession.



Notebooks and command-line tools automatically create a SparkSession for you



But you must create it yourself in all other environments.



We learned the bare minimum code for creating a Spark Session



I also talked about the session builder and the SparkContext



Thank You  
ScholarNest Technologies Pvt Ltd.  
[www.scholarnest.com](http://www.scholarnest.com)



ScholarNest

# Spark Azure Databricks

Databricks Spark Certification and beyond

Instructor: Prashant Kumar Pandey



Absolute Beginner to Specialization in Apache Spark and Azure Databricks





# Creating Data Frames Programmatically

We learned that Spark session is the starting point. However, Spark programming is all about reading, processing, and writing data. And Spark Dataframe is your tool for doing all of this.

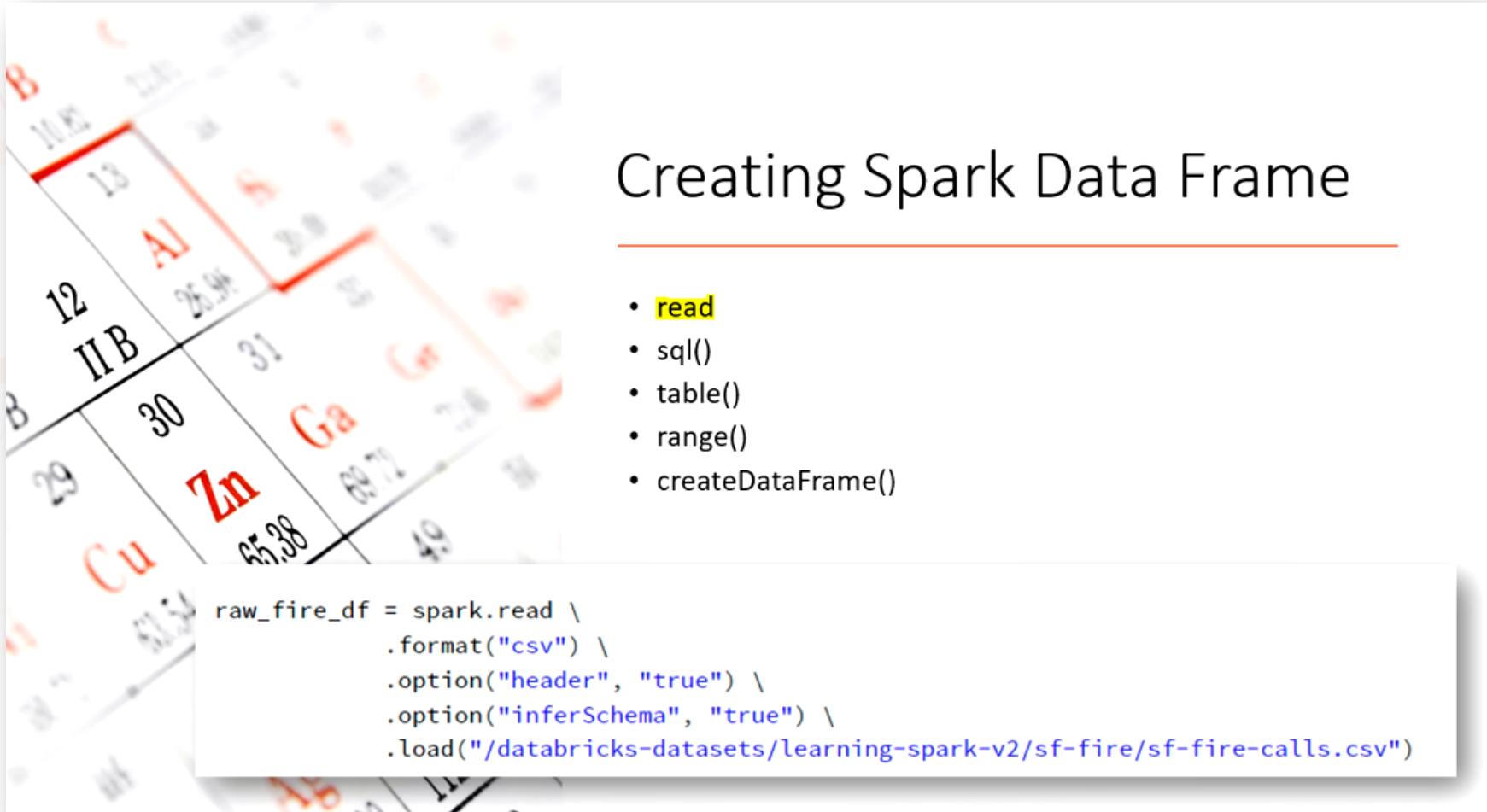
But how do you create a Spark Dataframe?

We have many ways to create a Spark Dataframe. Here is a list of some approaches:

1. `read`
2. `sql()`
3. `table()`
4. `range()`
5. `createDataFrame()`

All of these are offered by `SparkSession`.

We already learned about the Spark Session read attribute. Here is the code shown below that we used in an earlier example. The *SparkSession.read* gives you a *DataframeReader*, and you can use it to read data from a variety of data sources. We learned to read data from a CSV file, but *DataframeReader* supports many other sources.



## Creating Spark Data Frame

- **read**
- **sql()**
- **table()**
- **range()**
- **createDataFrame()**

```
raw_fire_df = spark.read \
    .format("csv") \
    .option("header", "true") \
    .option("inferSchema", "true") \
    .load("/databricks-datasets/learning-spark-v2/sf-fire/sf-fire-calls.csv")
```

We also learned to use the `SparkSession.sql()` method. Here is the earlier example. So you can use a Spark SQL in the `SparkSession.sql()` method, and the result of the SQL is a Spark Dataframe.

## Creating Spark Data Frame

- `read`
- `sql()`
- `table()`
- `range()`
- `createDataFrame()`

```
q1_sql_df = spark.sql("""  
    select count(distinct CallType) as distinct_call_type_count  
    from fire_service_calls_view  
    where CallType is not null  
    """)
```

The third method is even more straightforward than the SQL() method. You can use `spark.table()` method and pass in a table name to create a Dataframe. And we have an example code shown below. So you can give a database and table name to read a Spark table and create a Dataframe out of it. So if someone is asking you to convert a Spark table to a Spark Dataframe, you can use `spark.table()` method.

## Creating Spark Data Frame

- `read`
- `sql()`
- `table()`
- `range()`
- `createDataFrame()`

```
df1 = spark.table("spark_db.table_name")
```

The main focus of this document was to help you understand the `range()` and the `createDataFrame()` methods.

These two methods allow you to create a Dataframe without reading data from a source.

## Creating Spark Data Frame

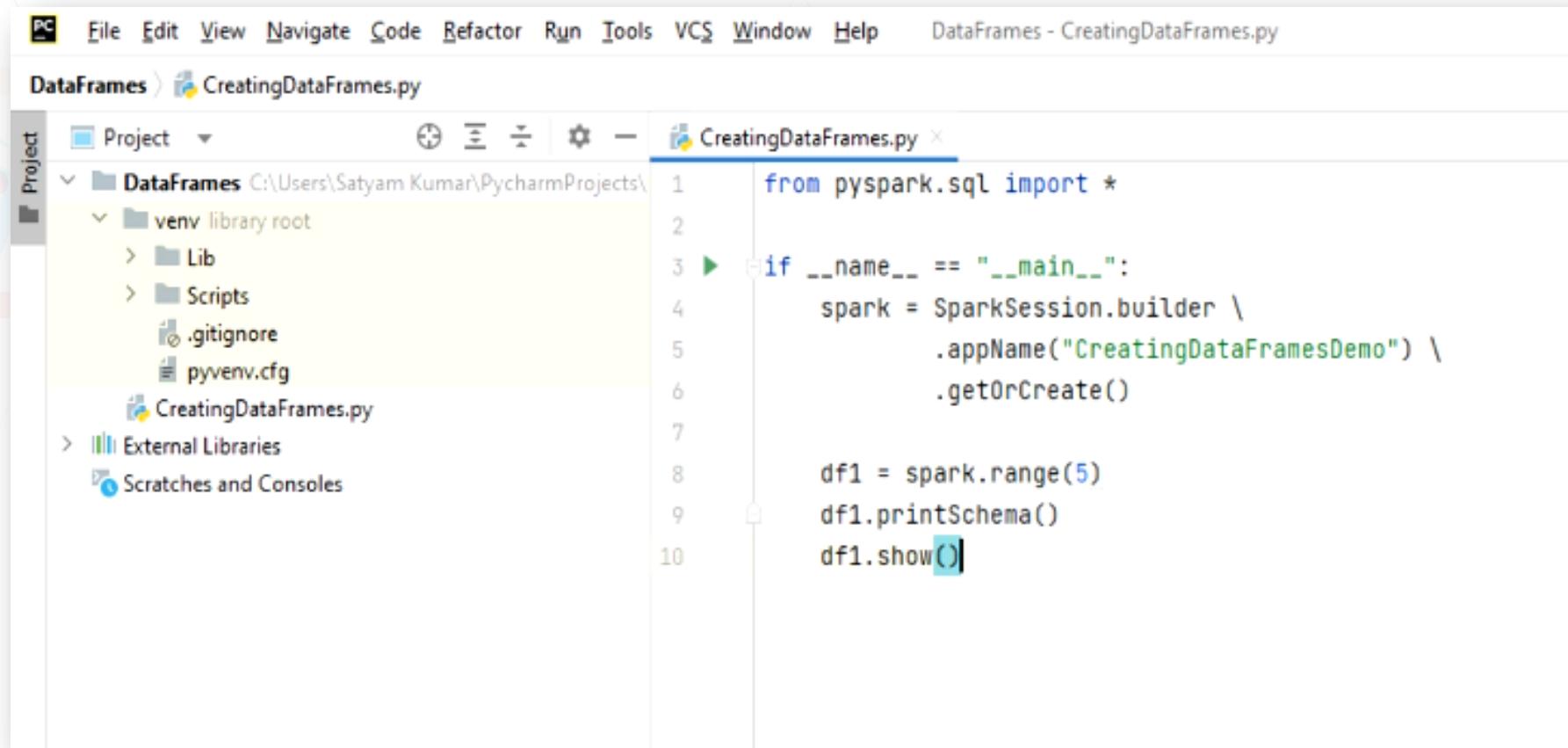
---

- `read`
- `sql()`
- `table()`
- `range()` ←
- `createDataFrame()` ←

Here I have opened a previous project which I created earlier. (**Reference – DataFrames**)

In this example, we learned to create a Spark Session. Now, the next step is to create a Data Frame.

I have written the code for that. We have a single-column Dataframe. The *spark.range()* allows you to create a single-column Dataframe. You might need a single-column data frame for testing something. And *spark.range()* is an easy method to create it. The range() method takes many parameters like start, end, step, and the number of partitions. You can show this Dataframe and check the schema.



The screenshot shows the PyCharm IDE interface. The top menu bar includes File, Edit, View, Navigate, Code, Refactor, Run, Tools, VCS, Window, Help, and the current file name, DataFrames - CreatingDataFrames.py. Below the menu is a toolbar with icons for Project, Run, Stop, and others. The left sidebar is the Project tool window, showing a 'DataFrames' project structure. Inside 'DataFrames', there is a 'venv' folder containing 'Lib', 'Scripts', '.gitignore', and 'pyenv.cfg'. The 'CreatingDataFrames.py' file is listed under 'Scripts'. Other sections in the Project window include 'External Libraries' and 'Scratches and Consoles'. The main editor window displays the following Python code:

```
from pyspark.sql import *
if __name__ == "__main__":
    spark = SparkSession.builder \
        .appName("CreatingDataFramesDemo") \
        .getOrCreate()

    df1 = spark.range(5)
    df1.printSchema()
    df1.show()
```

If we run the code, we can see the results as shown below.

By default, the column name is id, and its data type is long. The Dataframe values start from zero and stop at 5.

The screenshot shows the PyCharm IDE interface with the following details:

- File Menu:** PC, File, Edit, View, Navigate, Code, Refactor, Run, Tools, VCS, Window, Help.
- Project Tab:** DataFrames, CreatingDataFrames.py.
- Code Editor:** The file `CreatingDataFrames.py` contains the following Python code:

```
df1 = spark.range(5)
df1.printSchema()
df1.show()

if __name__ == "__main__":
    pass
```
- Run Tab:** The output of the run command is displayed:

```
22/04/24 13:07:53 WARN ProcfsMetricsGetter: Exception when trying to compute pagesize, as a result reporting of ProcessTree metrics is stopped
root
|-- id: long (nullable = false)

+---+
| id|
+---+
| 0|
| 1|
| 2|
| 3|
| 4|
+---+
```
- Bottom Status:** Process finished with exit code 0.

If you find it difficult to work with IDE, you can switch to Notebooks.

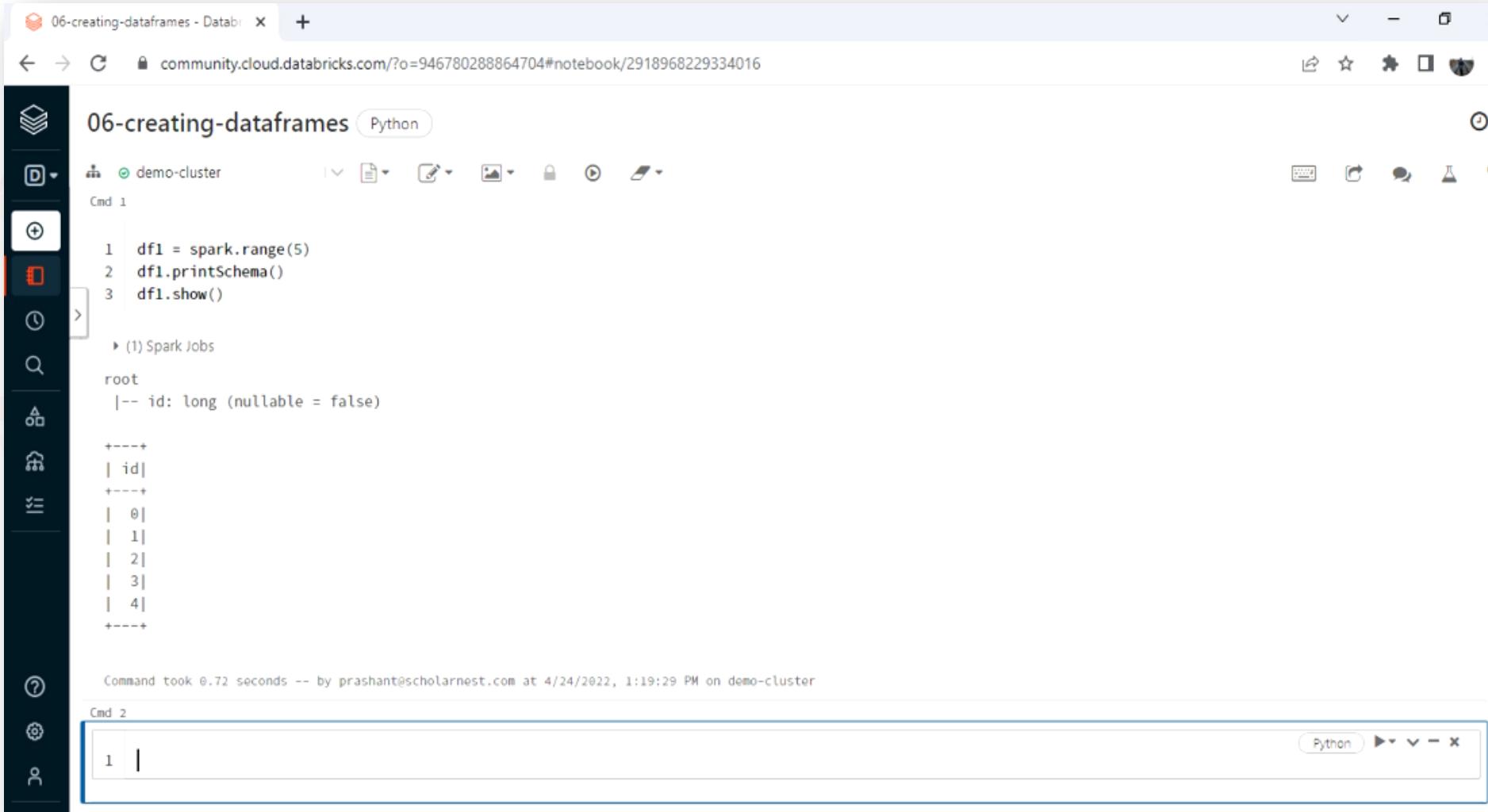
Here I am in my Databricks Workspace.

Create a new notebook. (**Reference - 06-creating-dataframes.ipynb**)

The screenshot shows the Databricks Community Edition workspace. The left sidebar contains navigation icons for notebooks, clusters, files, and more. The main area is titled "Data Science & Engineering". It features three main sections: "Notebook" (with a "Create a notebook" button), "Data import" (with a "Browse files" button), and "Partner Connect" (listing Fivetran, dbt, Tableau, and Power BI). Below these are "Guide: Quickstart tutorial" and "Recents". The "Recents" section lists five notebooks with their names and last viewed times:

Name	Last viewed
05-working-with-dataframe	16 days ago
03-spark-table-demo	16 days ago
04-spark-sql-demo	16 days ago
02-spark-dataframe-demo	17 days ago
01-getting-started	17 days ago

Now you can go back to the IDE and copy the last three lines of the code. Databricks notebook automatically creates a Spark session variable for me. So I do not need to create a Spark session variable. If you run the cell, you will get the same output.



The screenshot shows a Databricks notebook titled "06-creating-dataframes" in Python. The notebook URL is [community.cloud.databricks.com/?o=946780288864704#notebook/2918968229334016](https://community.cloud.databricks.com/?o=946780288864704#notebook/2918968229334016). The sidebar on the left contains icons for file operations, clusters, search, and help. The main area shows a command cell (Cmd 1) with the following Python code:

```
1 df1 = spark.range(5)
2 df1.printSchema()
3 df1.show()
```

When run, the cell produces the following output:

```
+-----+
| id|
+---+
| 0|
| 1|
| 2|
| 3|
| 4|
+---+
```

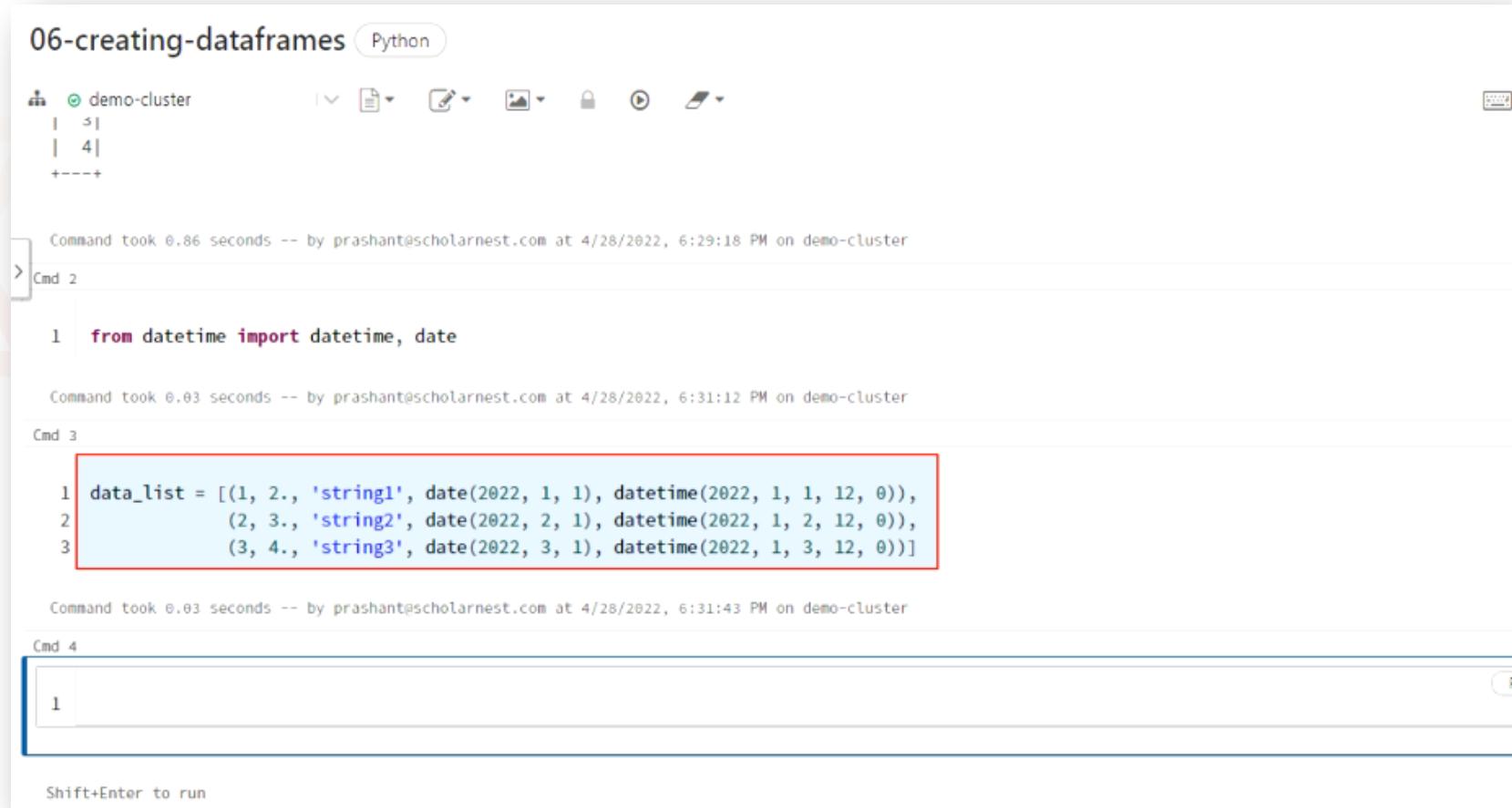
A message at the bottom of the cell indicates: "Command took 0.72 seconds -- by prashant@scholarnest.com at 4/24/2022, 1:19:29 PM on demo-cluster". Below this is another command cell (Cmd 2) which is currently empty.

The range method is basic. It only supports creating a single-column Dataframe. But we also have `spark.createDataFrame()` method. It converts a list into a Spark Dataframe. Look at the requirement given in the image below to create a Dataframe programmatically. You have column names shown here. Then you have some rows shown in the diagram. You also expect column data types, as shown at the top.

long	double	string	date	timestamp
a	b	c	d	e
1	2.0	string1	2000-01-01	2000-01-01 12:00:00
2	3.0	string2	2000-02-01	2000-01-02 12:00:00
3	4.0	string3	2000-03-01	2000-01-03 12:00:00

The `createDataFrame()` takes a Python list and converts it into a Dataframe. We define a python list inside a squire brackets. Then we define each record as a tuple inside a pair of parenthesis. I wanted three records, so I defined three types inside the squire bracket.

The column is a date, and the last one is a datetime, so I am using Python date and datetime functions from the Python datetime module. You must import the date and datetime functions from the datetime module.



The screenshot shows a Jupyter Notebook interface with the title "06-creating-dataframes" and a Python tab selected. The sidebar shows a "demo-cluster" with two files: "5" and "4".

Cmd 2:

```
1 from datetime import datetime, date
```

Command took 0.03 seconds -- by prashant@scholarnest.com at 4/28/2022, 6:31:12 PM on demo-cluster

Cmd 3:

```
1 data_list = [(1, 2., 'string1', date(2022, 1, 1), datetime(2022, 1, 1, 12, 0)),
 2             (2, 3., 'string2', date(2022, 2, 1), datetime(2022, 1, 2, 12, 0)),
 3             (3, 4., 'string3', date(2022, 3, 1), datetime(2022, 1, 3, 12, 0))]
```

1 data\_list = [(1, 2., 'string1', date(2022, 1, 1), datetime(2022, 1, 1, 12, 0)),  
2 (2, 3., 'string2', date(2022, 2, 1), datetime(2022, 1, 2, 12, 0)),  
3 (3, 4., 'string3', date(2022, 3, 1), datetime(2022, 1, 3, 12, 0))]

Command took 0.03 seconds -- by prashant@scholarnest.com at 4/28/2022, 6:31:43 PM on demo-cluster

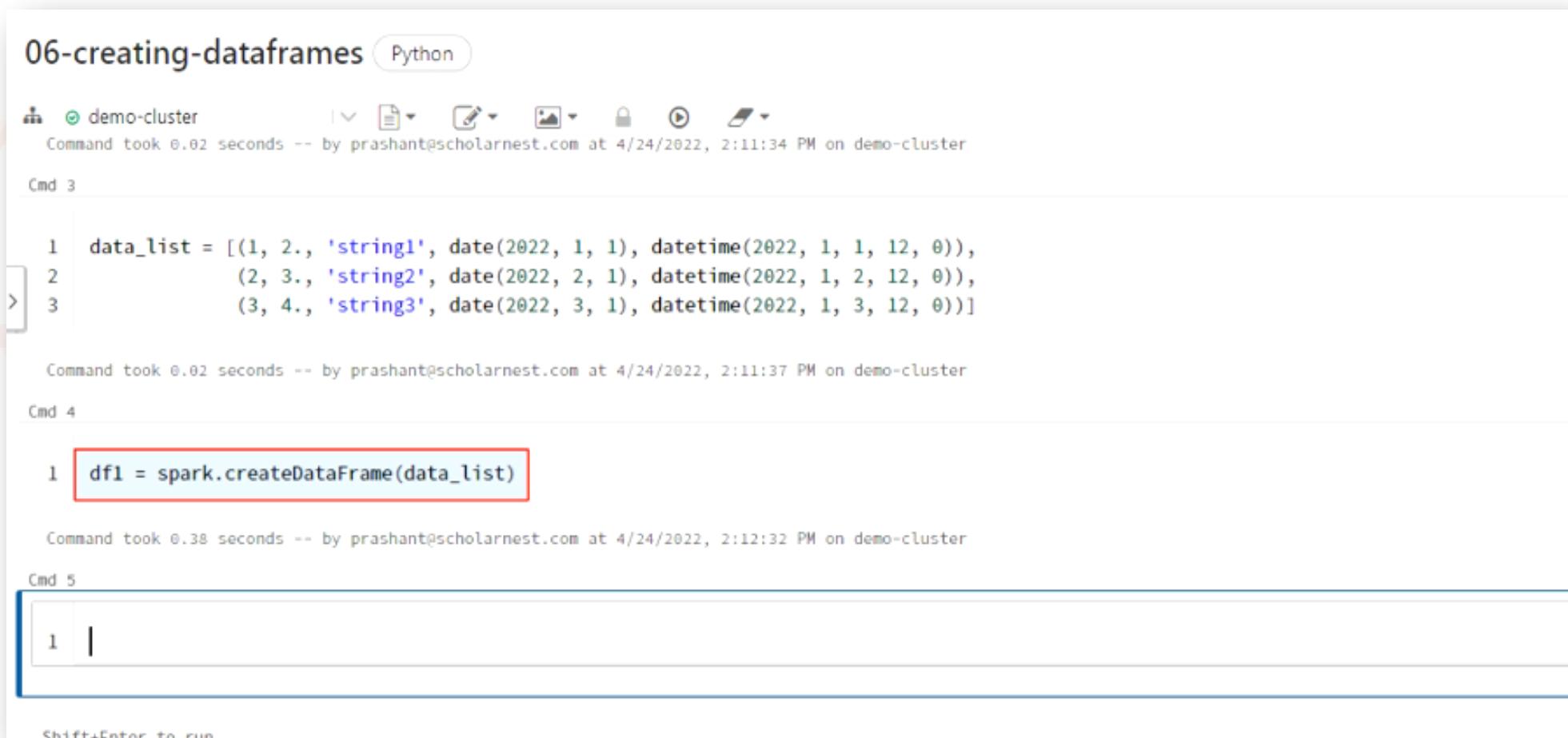
Cmd 4:

```
1
```

Shift+Enter to run

Once you have the list, you can create a Dataframe as shown below. I am using `spark.createDataFrame()` method and passing the `data_list`.

The `createDataFrame()` method will convert the list into a Spark Dataframe and return it. So the `df1` is my Spark Dataframe.



06-creating-dataframes Python

demo-cluster Command took 0.02 seconds -- by prashant@scholarnest.com at 4/24/2022, 2:11:34 PM on demo-cluster

Cmd 3

```
1 data_list = [(1, 2., 'string1', date(2022, 1, 1), datetime(2022, 1, 1, 12, 0)),  
2                 (2, 3., 'string2', date(2022, 2, 1), datetime(2022, 1, 2, 12, 0)),  
3                 (3, 4., 'string3', date(2022, 3, 1), datetime(2022, 1, 3, 12, 0))]
```

Command took 0.02 seconds -- by prashant@scholarnest.com at 4/24/2022, 2:11:37 PM on demo-cluster

Cmd 4

```
1 df1 = spark.createDataFrame(data_list)
```

Command took 0.38 seconds -- by prashant@scholarnest.com at 4/24/2022, 2:12:32 PM on demo-cluster

Cmd 5

```
1 |
```

Shift+Enter to run

You can show the Dataframe and check if you got the desired output. Adjacent to the output is an image of the desired Dataframe. I am almost there except for a tiny problem. The column names are not as per my requirement. It took some default column names here. But I wanted to give some different column names. I will fix this in a while, but before that let us check the schema.

06-creating-dataframes Python

demo-cluster

Cmd 4

```
1 df1 = spark.createDataFrame
```

Command took 0.38 seconds -- by prashant@scholarnest.com

Cmd 5

```
1 df1.show()
```

(3) Spark Jobs

long	double	string	date	timestamp
a	b	c	d	e
1	2.0	string1	2000-01-01	2000-01-01 12:00:00
2	3.0	string2	2000-02-01	2000-01-02 12:00:00
3	4.0	string3	2000-03-01	2000-01-03 12:00:00

_1	_2	_3	_4	_5
1	2.0	string1	2022-01-01	2022-01-01 12:00:00
2	3.0	string2	2022-02-01	2022-01-02 12:00:00
3	4.0	string3	2022-03-01	2022-01-03 12:00:00

Command took 1.28 seconds -- by prashant@scholarnest.com at 4/24/2022, 2:12:55 PM on demo-cluster

Here is the schema shown below. It looks like data types are fine. The *createDataFrame* automatically infers the schema of the Dataframe based on the given data. So do remember to define your data list with correct data types so the *createDataFrame* can easily infer it.

The screenshot shows a Jupyter Notebook cell with the title "06-creating-dataframes" and the language "Python". The cell contains the following code:

```
1 df1.printSchema()
```

Below the code, the output shows the schema of the DataFrame:

```
root
|-- _1: long (nullable = true)
|-- _2: double (nullable = true)
|-- _3: string (nullable = true)
|-- _4: date (nullable = true)
|-- _5: timestamp (nullable = true)
```

The last three columns (\_3, \_4, \_5) are highlighted with a red box.

At the bottom of the cell, the command took 0.01 seconds and was run by prashant@scholarnest.com at 4/24/2022, 2:13:38 PM on demo-cluster.

Now let's come back to the column name problem. How do I fix these column names? One easy way is to use the *toDF()* method. The *toDF()* is a Dataframe transformation. So once you have the Dataframe, you can use the *toDF()* to change all the column names at once. You can chain it after the *createDataFrame()* as shown below.

The screenshot shows a Jupyter Notebook cell titled "06-creating-dataframes" in Python mode. The cell contains the following code:

```
1 from datetime import datetime, time
```

Command took 0.02 seconds -- by prashant@scholarnest.com at 4/24/2022, 2:11:34 PM on demo-cluster

```
1 data_list = [(1, 2., 'string1', date(2022, 1, 1), datetime(2022, 1, 1, 12, 0)),  
2                 (2, 3., 'string2', date(2022, 2, 1), datetime(2022, 1, 2, 12, 0)),  
3                 (3, 4., 'string3', date(2022, 3, 1), datetime(2022, 1, 3, 12, 0))]
```

Command took 0.02 seconds -- by prashant@scholarnest.com at 4/24/2022, 2:11:37 PM on demo-cluster

```
1 df1 = spark.createDataFrame(data_list).toDF("a", "b", "c", "d", "e")|
```

Command took 0.38 seconds -- by prashant@scholarnest.com at 4/24/2022, 2:12:32 PM on demo-cluster

The line `df1 = spark.createDataFrame(data_list).toDF("a", "b", "c", "d", "e")` is highlighted with a red border.

You can again check your Dataframe, and the column names are rectified as desired.

06-creating-dataframes Python

demo-cluster Command took 0.09 seconds -- by prashant@scholarnest.com at 4/24/2022, 2:15:17 PM on demo-cluster

Cmd 5

```
1 df1.show()
```

▶ (3) Spark Jobs

a	b	c	d	e
1	2.0	string1	2022-01-01	2022-01-01 12:00:00
2	3.0	string2	2022-02-01	2022-01-02 12:00:00
3	4.0	string3	2022-03-01	2022-01-03 12:00:00

Command took 0.52 seconds -- by prashant@scholarnest.com at 4/24/2022, 2:15:21 PM on demo-cluster

Cmd 6

We created a *data\_list* and used it to create a Dataframe. The *createDataFrame()* automatically infers the schema and generates some default column names. But you can supply a schema definition and enforce it. And here are two ways to define a schema as shown below.

The first one is simply a list of column names.

The second approach defines the column names and data types both.

The screenshot shows a Jupyter Notebook cell with the title "06-creating-dataframes" and the language "Python". The cell displays two schema definitions:

```
schema_1 = ['a', 'b', 'c', 'd', 'e']
schema_2 = 'a int, b double, c string, d date, e timestamp'
```

Below the code, the output shows the schema inferred by the system:

```
|--- a: string (nullable = true)
|--- b: date (nullable = true)
|--- c: timestamp (nullable = true)
```

At the bottom, the command took 0.04 seconds to run.

I am using the first schema to show you the result.

I supplied a list of column names, so the *createDataFrame* should use the given columns names instead of the default names. Data types are still correct because the *createDataFrame()* could infer it correctly from the *data\_list*.

06-creating-dataframes Python

demo-cluster    Command took 0.04 seconds -- by prashant@scholarnest.com at 4/24/2022, 2:19:38 PM on demo-cluster

Cmd: 7

```
1 schema_1 = ['a', 'b', 'c', 'd', 'e']
2 schema_2 = 'a int, b double, c string, d date, e timestamp'
```

Command took 0.04 seconds -- by prashant@scholarnest.com at 4/24/2022, 2:19:38 PM on demo-cluster

Cmd: 8

```
1 spark.createDataFrame(data_list, schema_1).printSchema()
```

root

```
|-- a: long (nullable = true)
|-- b: double (nullable = true)
|-- c: string (nullable = true)
|-- d: date (nullable = true)
|-- e: timestamp (nullable = true)
```

Command took 0.09 seconds -- by prashant@scholarnest.com at 4/24/2022, 2:20:28 PM on demo-cluster

Then in the second schema, you can see that I purposefully change the data type for the first column. I know the `createDataFrame` will infer the first column to be a long value. But I want it to be an integer. I tried enforcing a schema and stopped `createDataFrame` from inferring it, and it worked. You can see the output in the bottom section of the image below.



```
06-creating-dataframes Python

demo-cluster Cmd 8

1 spark.createDataFrame(data_list, schema_1).printSchema()

root
 |-- a: long (nullable = true)
 |-- b: double (nullable = true)
 |-- c: string (nullable = true)
 |-- d: date (nullable = true)
 |-- e: timestamp (nullable = true)

Command took 0.09 seconds -- by prashant@scholarnest.com at 4/24/2022, 2:20:28 PM on demo-cluster

Cmd 9

1 spark.createDataFrame(data_list, schema_2).printSchema()

root
 |-- a: integer (nullable = true)
 |-- b: double (nullable = true)
 |-- c: string (nullable = true)
 |-- d: date (nullable = true)
 |-- e: timestamp (nullable = true)

Command took 0.15 seconds -- by prashant@scholarnest.com at 4/24/2022, 2:21:32 PM on demo-cluster
```

The *createDataFrame* is even more flexible.

We used Python list and schema definition for creating a Spark Dataframe.

However, the create Dataframe can also take the following.

1. List of DataFrame Row
2. Pandas DataFrame
3. Spark RDD

## List of DataFrame Row:

I am importing the Row class from the pyspark.sql package. You can create a list of rows and use them to create a Dataframe programmatically. Each Row is a comma-separated list of column names and column value pairs. Once you have the row list, you can use it in the *createDataFrame()* method as we used the python list.



The screenshot shows a Jupyter Notebook cell with the title "06-creating-dataframes" and the language selector "Python". The cell contains the following code:

```
1 from pyspark.sql import Row
2
3 row_list = [Row(a=1, b=2.0, c='string1', d=date(2022, 1, 1), e=datetime(2022, 1, 1, 12, 0)),
4             Row(a=2, b=3.0, c='string2', d=date(2022, 2, 1), e=datetime(2022, 1, 2, 12, 0)),
5             Row(a=3, b=4.0, c='string3', d=date(2022, 3, 1), e=datetime(2022, 1, 3, 12, 0))]
```

After running the cell, the output shows the command took 0.15 seconds and was run by prashant@scholarnest.com at 4/24/2022, 2:21:32 PM on demo-cluster. The cell is numbered "Cmd 10".

## List of DataFrame Row:

You can print the schema to see that row list comes with the column names, so the *createDataFrame()* method takes the correct names. However, it will infer the schema. But you can also pass the schema to the *createDataFrame()* and enforce the data types.

The screenshot shows a Jupyter Notebook cell with the title "06-creating-dataframes" and the language "Python". The cell contains the following code:

```
1 from pyspark.sql import Row
2
3 row_list = [Row(a=1, b=2.0, c='string1', d=date(2022, 1, 1), e=datetime(2022, 1, 1, 12, 0)),
4             Row(a=2, b=3.0, c='string2', d=date(2022, 2, 1), e=datetime(2022, 1, 2, 12, 0)),
5             Row(a=3, b=4.0, c='string3', d=date(2022, 3, 1), e=datetime(2022, 1, 3, 12, 0))]
```

Below the code, the output shows the command took 0.03 seconds and was run by prashant@scholarnest.com at 4/24/2022, 2:26:04 PM on demo-cluster. The next command, "Cmd 11", is shown as:

```
1 spark.createDataFrame(row_list).printSchema()
```

The output of this command is:

```
root
 |-- a: long (nullable = true)
 |-- b: double (nullable = true)
 |-- c: string (nullable = true)
 |-- d: date (nullable = true)
 |-- e: timestamp (nullable = true)
```

Below this, the command took 0.08 seconds and was run by prashant@scholarnest.com at 4/24/2022, 2:26:34 PM on demo-cluster.

## Pandas DataFrame:

So I am importing the panda's package and defining a panda Dataframe. You might not have worked with pandas package, and that's fine. We are just showing an example, so you know how to convert a pandas Dataframe into a Spark Dataframe.

The screenshot shows a Jupyter Notebook cell with the title "06-creating-dataframes". The cell contains the following Python code:

```
1 import pandas as pd
2
3 pd_df = pd.DataFrame({'a':[1, 2, 3],
4                       'b':[2.0, 3.0, 4.0],
5                       'c':['string1', 'string2', 'string3'],
6                       'd':[date(2022, 1, 1), date(2022, 2, 1), date(2022, 3, 1)],
7                       'e':[datetime(2022, 1, 1, 12, 0), datetime(2022, 1, 2, 12, 0), datetime(2022, 1, 3, 12, 0)]})
```

After running the code, the output shows the schema of the DataFrame:

```
|--- c: string (nullable = true)
|--- d: date (nullable = true)
|--- e: timestamp (nullable = true)
```

Execution details are shown at the bottom:

```
Command took 0.08 seconds -- by prashant@scholarnest.com at 4/24/2022, 2:26:34 PM on demo-cluster
> Cmd 12
```

Command took 0.03 seconds -- by prashant@scholarnest.com at 4/24/2022, 2:30:37 PM on demo-cluster

## Pandas DataFrame:

You can use the `createDataFrame()` to print the schema. I used `schema_2` so I could enforce the schema on this Dataframe. However, the schema is an optional argument. If you do not supply the schema, the spark will infer it from the given values.

```
06-creating-dataframes Python

demo-cluster

1 import pandas as pd
2
3 pd_df = pd.DataFrame({'a':[1, 2, 3],
4                      'b':[2.0, 3.0, 4.0],
5                      'c':['string1', 'string2', 'string3'],
6                      'd':[date(2022, 1, 1), date(2022, 2, 1), date(2022, 3, 1)],
7                      'e':[datetime(2022, 1, 1, 12, 0), datetime(2022, 1, 2, 12, 0), datetime(2022, 1, 3, 12, 0)]})

Command took 0.03 seconds -- by prashant@scholarnest.com at 4/24/2022, 2:30:37 PM on demo-cluster
Cmd 13

1 spark.createDataFrame(pd_df, schema_2).printSchema()

root
|-- a: integer (nullable = true)
|-- b: double (nullable = true)
|-- c: string (nullable = true)
|-- d: date (nullable = true)
|-- e: timestamp (nullable = true)

Command took 0.32 seconds -- by prashant@scholarnest.com at 4/24/2022, 2:31:26 PM on demo-cluster
```

## **Spark RDD:**

Spark core package defines an RDD data structure. The full form of the RDD is Resilient Distributed Dataset. However, it is as simple as an in-memory distributed data structure. In the earlier versions of Spark, we used RDDs for data processing. However, the later versions of Spark came up with the Dataframe. Dataframe is also a distributed data structure. It is built on top of the RDD. However, Dataframe is optimized and more efficient than the RDD.

The RDD cannot have a schema. It is schema-less. But Dataframe must have a schema. You cannot create a schema-less Dataframe. It must have a schema. If you do not provide a schema, Spark will apply some default schema. We saw the same earlier, *createDataFrame()* method infers the data types and adds some default column names to create a default schema.

Hence, Spark Dataframe is a more advanced, optimized, and schema-full data structure. We recommend using Dataframes as much as possible. Spark RDD is part of the Spark core package, and it is schema-less. We should avoid using RDD and use Dataframes.

## Spark RDD:

Here is the code for manually creating an RDD. I am using the spark session to get the SparkContext. Why SparkContext? Because SparkContext is my gateway to the Spark core package. Then we use the parallelize() method to create an RDD.

Once you have an RDD, you can impose a schema and convert it to a Dataframe.

06-creating-dataframes Python

```
demo-cluster
|-- c: string (nullable = true)
|-- d: date (nullable = true)
|-- e: timestamp (nullable = true)

Command took 0.32 seconds -- by prashant@scholarnest.com at 4/24/2022, 2:31:26 PM on demo-cluster
> Cmd 14

1 rdd = spark.sparkContext.parallelize([(1, 2., 'string1', date(2022, 1, 1), datetime(2022, 1, 1, 12, 0)),
2                                         (2, 3., 'string2', date(2022, 2, 1), datetime(2022, 1, 2, 12, 0)),
3                                         (3, 4., 'string3', date(2022, 3, 1), datetime(2022, 1, 3, 12, 0))])

Command took 0.05 seconds -- by prashant@scholarnest.com at 4/24/2022, 2:34:57 PM on demo-cluster
```

# Spark RDD:

You can see the desired output if you print the schema.



06-creating-dataframes Python

demo-cluster

Cmd 14

```
1 rdd = spark.sparkContext.parallelize([(1, 2., 'string1', date(2022, 1, 1), datetime(2022, 1, 1, 12, 0)),
2                                         (2, 3., 'string2', date(2022, 2, 1), datetime(2022, 1, 2, 12, 0)),
3                                         (3, 4., 'string3', date(2022, 3, 1), datetime(2022, 1, 3, 12, 0))])
```

Command took 0.05 seconds -- by prashant@scholarnest.com at 4/24/2022, 2:34:57 PM on demo-cluster

Cmd 15

```
1 spark.createDataFrame(rdd, schema_1).printSchema()
```

▶ (2) Spark Jobs

```
root
|-- a: long (nullable = true)
|-- b: double (nullable = true)
|-- c: string (nullable = true)
|-- d: date (nullable = true)
|-- e: timestamp (nullable = true)
```

Command took 0.70 seconds -- by prashant@scholarnest.com at 4/24/2022, 2:36:21 PM on demo-cluster



Thank You  
ScholarNest Technologies Pvt Ltd.  
[www.scholarnest.com](http://www.scholarnest.com)



ScholarNest

# Spark Azure Databricks

Databricks Spark Certification and beyond

Instructor: Prashant Kumar Pandey



Absolute Beginner to Specialization in Apache Spark and Azure Databricks





# Spark Dataframes Schema

## Requirement:

You are given a flight time data set. This data set comes in a JSON file, and you are asked to load this Dataframe and create a Spark table. For your reference, here is one sample record, and here are details about the columns shown below.

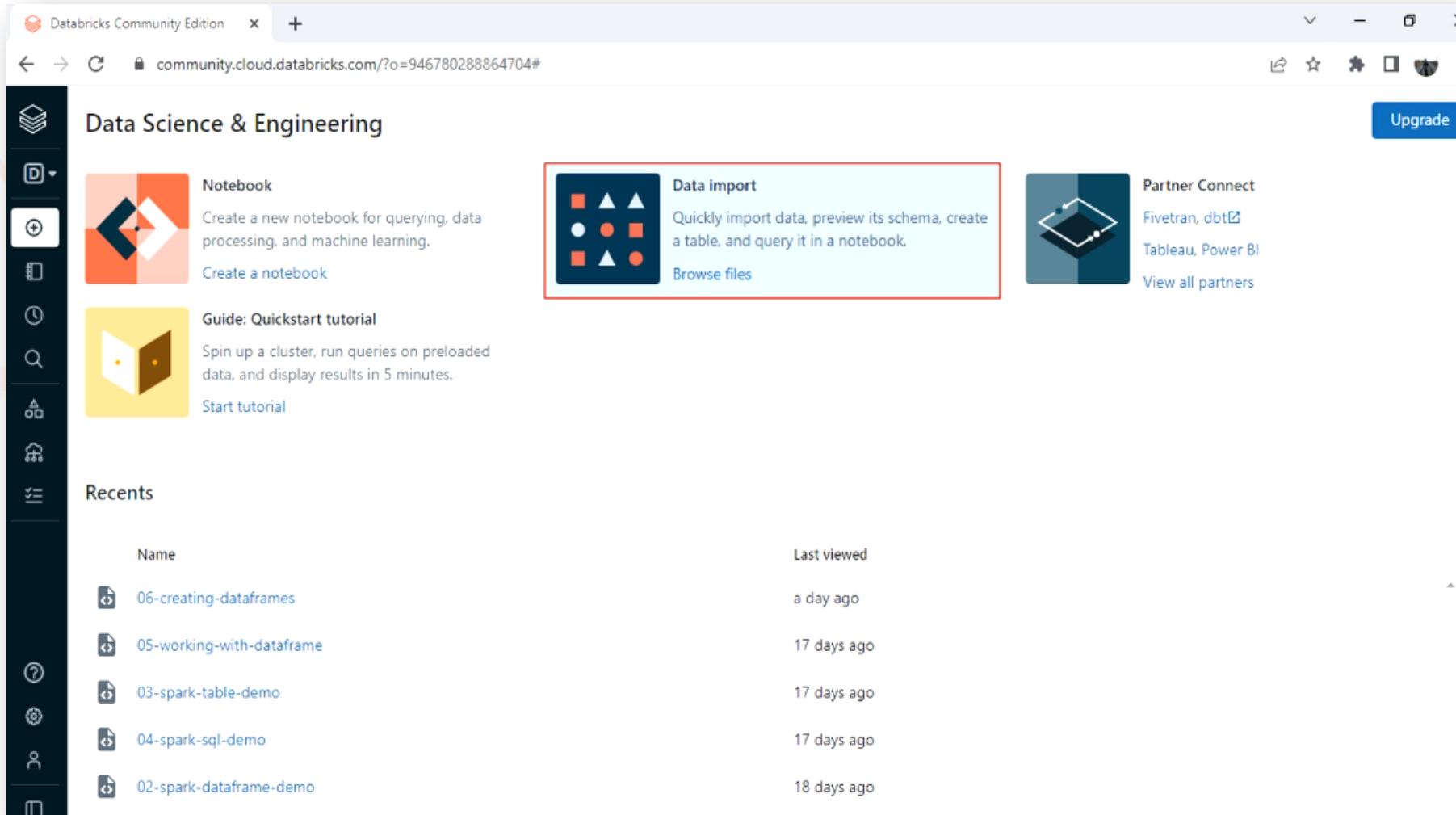
### Flight Delay Data

```
{  
    "FL_DATE": "1/1/2000",  
    "OP_CARRIER": "DL",  
    "OP_CARRIER_FL_NUM": 1451,  
    "ORIGIN": "BOS",  
    "ORIGIN_CITY_NAME": "Boston, MA",  
    "DEST": "ATL",  
    "DEST_CITY_NAME": "Atlanta, GA",  
    "CRS_DEP_TIME": 1115,  
    "DEP_TIME": 1113,  
    "WHEELS_ON": 1343,  
    "TAXI_IN": 5,  
    "CRS_ARR_TIME": 1400,  
    "ARR_TIME": 1348,  
    "CANCELLED": 0,  
    "DISTANCE": 946  
}
```

1. Flight date in MM/DD/YYYY format
2. Unique carrier code
3. Flight number
4. Origin IATA Airport Code
5. Origin City Name
6. Destination IATA Airport Code
7. Destination City Name
8. Scheduled Departure time in HHMM
9. Actual Departure time in HHMM
10. Flight Duration in HHMM
11. TaxiIn taxi in time, in minutes
12. Scheduled Arrival time in HHMM
13. Actual Departure time in HHMM
14. Was it canceled? 1 for a true, and 0 is false
15. Distance in miles

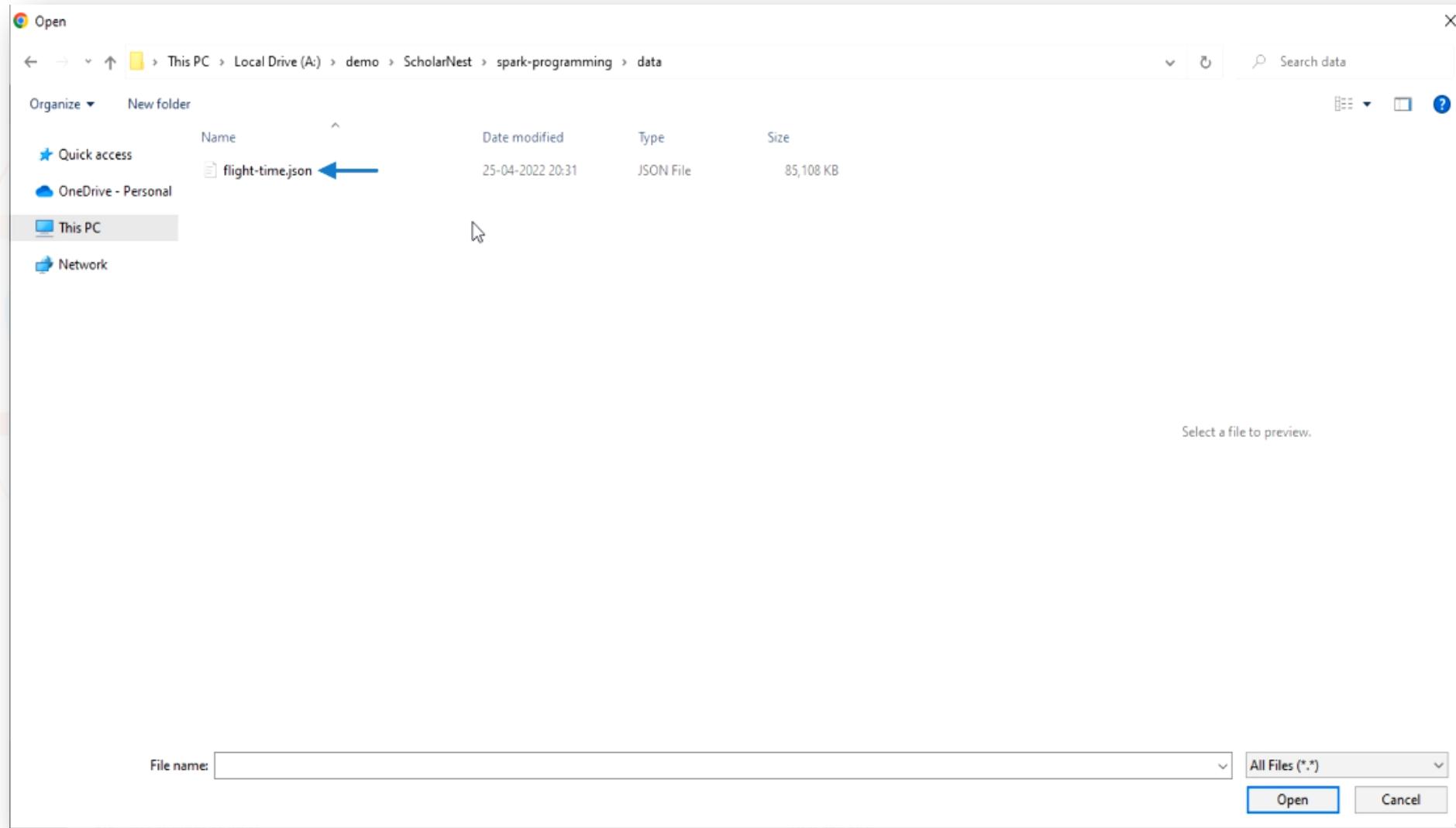
# Load the table into a Dataframe:

Go to your Databricks workspace. My data file is available on my local machine. I will upload the data file to the cloud so I can work with it. Click this Data import button highlighted below.



# Load the table into a Dataframe:

Navigate to the folder that contains the file and upload it.



## Load the table into a Dataframe:

Your file is uploaded. It goes into the Databricks File System and sits in the `/FileStore/tables/` directory. The Databricks File System is popularly known as DBFS, and the community edition DBFS is built on top of the Amazon S3.

### Create New Table

Data source [?](#)

Upload File   S3   DBFS   Other Data Sources   Partner Integrations

DBFS Target Directory [?](#)

/FileStore/tables/ (optional) [Select](#)

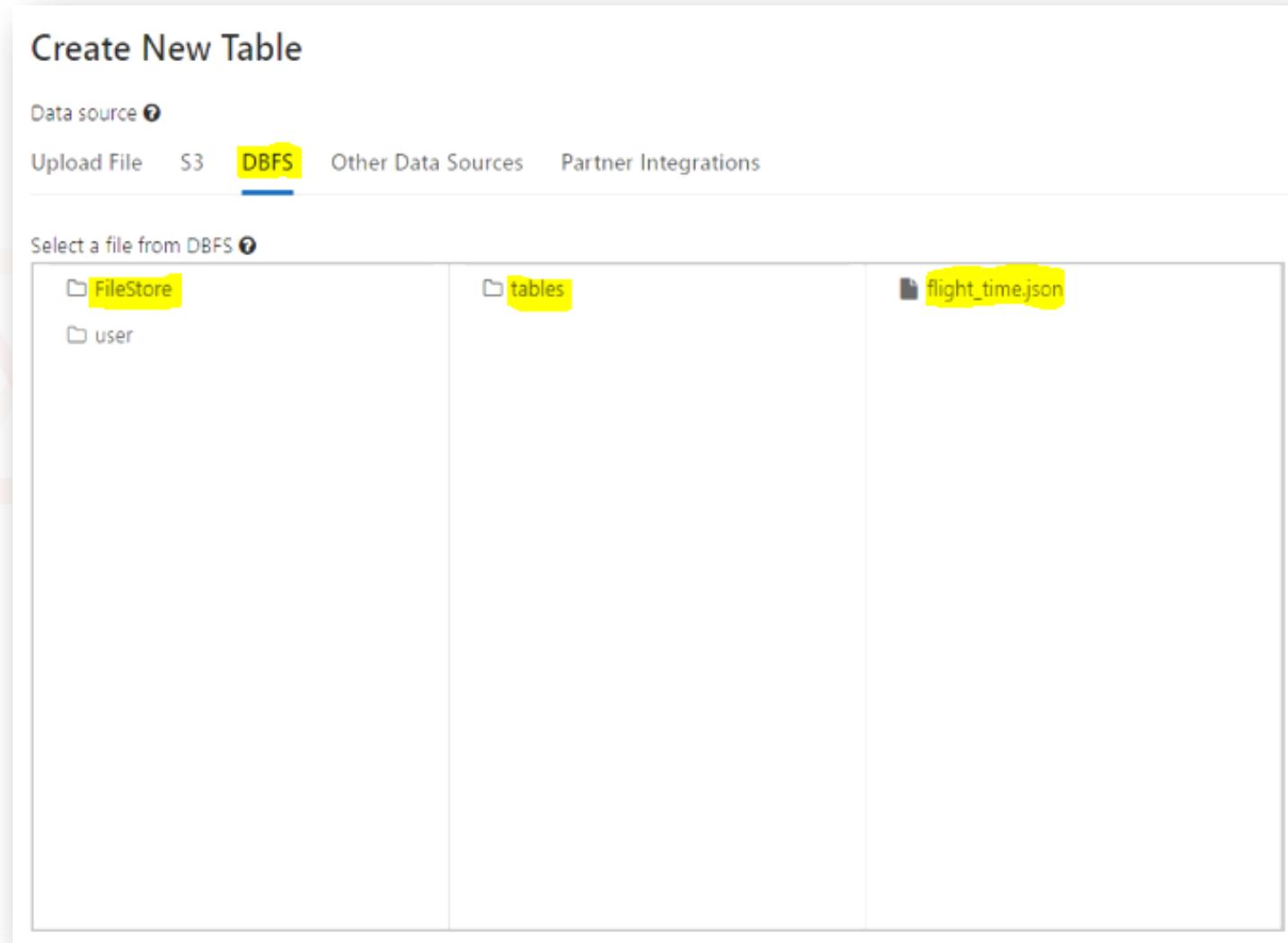
Files [?](#)

flight-time.json ✓  
87.2 MB [Remove file](#)

✓ File uploaded to /FileStore/tables/flight\_time.json 

[Create Table with UI](#)  [Create Table in Notebook](#) [?](#)

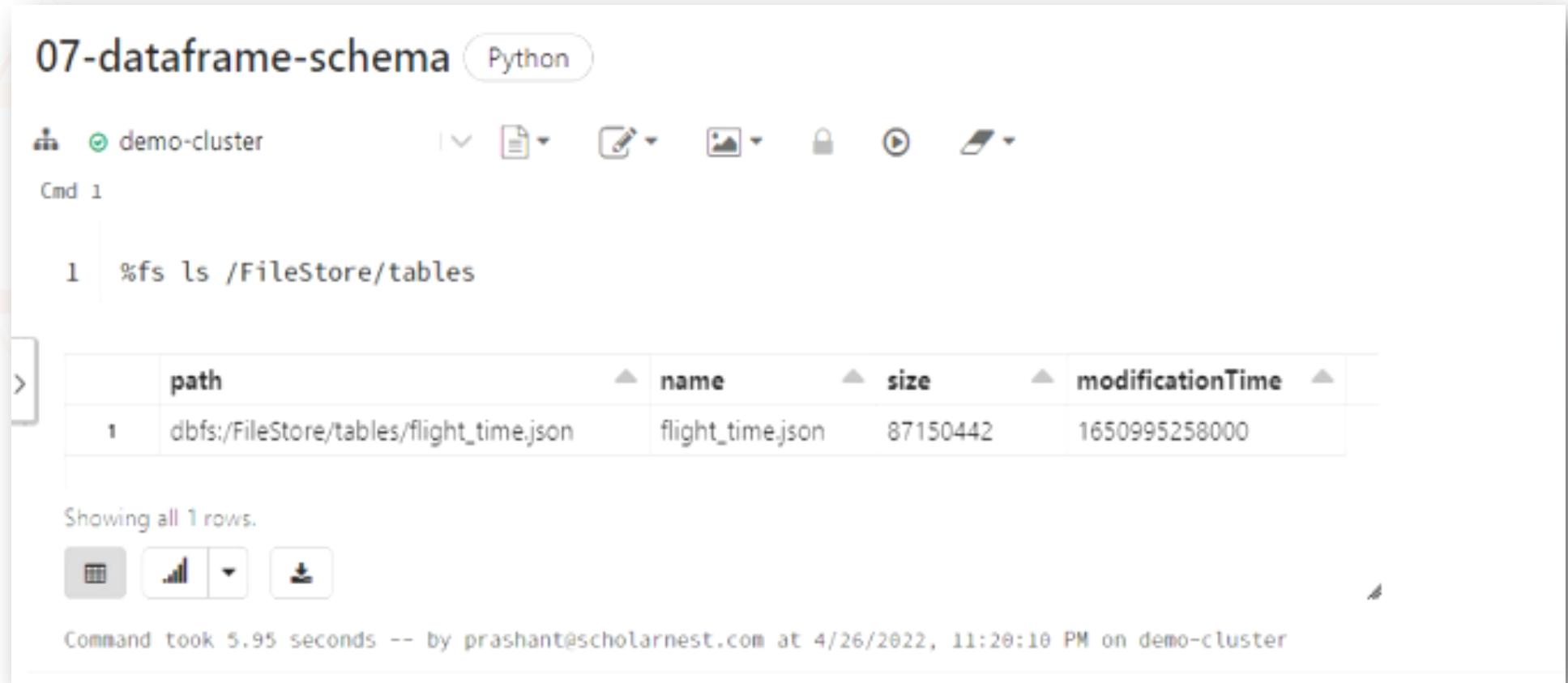
To check your file, go to the DBFS tab and browse to the *FileStore/tables/*, and you will see the file. So my data file is now available in the cloud. Now you can create a new notebook to write code to meet the given requirement.



## Create a new notebook. (**Reference - 07-dataframe-schema.ipynb**)

You can use the %fs ls command to check your data file, and we see the same in the output below.

The %fs ls is a DBFS command, and it is similar to the Linux ls command.



The screenshot shows a Jupyter Notebook interface with the title "07-dataframe-schema" and a Python tab selected. The toolbar includes icons for file operations like New, Open, Save, and Run. Below the toolbar, the cluster "demo-cluster" is selected. A command cell labeled "Cmd 1" contains the command "%fs ls /FileStore/tables". The output cell displays a table with one row, showing the file "flight\_time.json" located at "dbfs:/FileStore/tables/flight\_time.json". The table has columns: path, name, size, and modificationTime. The table shows the following data:

	path	name	size	modificationTime
1	dbfs:/FileStore/tables/flight_time.json	flight_time.json	87150442	1650995258000

Below the table, it says "Showing all 1 rows." and provides icons for cell operations like Cut, Copy, Paste, and Delete. At the bottom of the output cell, a message indicates the command took 5.95 seconds and was run by prashant@scholarnest.com at 4/26/2022, 11:20:10 PM on demo-cluster.

You have seen a similar code to create a Dataframe earlier.  
I changed the format from CSV to JSON. Because this time, I am loading JSON data.  
I also removed the header and inferSchema options. Because a JSON file doesn't come with a header row. The header option is only applicable for the CSV files.  
The inferSchema option is not applicable for JSON. Because a JSON document comes with the column names in every record. Spark will always read the column names. In fact, it will read the column name with every record.

Cmd 2

```
1 flight_time_raw_df = spark.read \
2           .format("json") \
3           .load("/FileStore/tables/flight_time.json")
```

▶ (1) Spark Jobs

Command took 20.42 seconds -- by prashant@scholarnest.com at 4/26/2022, 11:22:50 PM on demo-cluster

Look at the schema of the data frame we created. You can see that Spark automatically inferred the schema for this Dataframe. Some columns are long, and others are string. And that is not correct.

07-dataframe-schema Python

demo-cluster (1) Spark Jobs

Command took 20.42 seconds -- by prashant@scholarnest.com at 4/26/2022, 11:22:50 PM on demo-cluster

Cmd 3

```
> 1 flight_time_raw_df.printSchema()

root
|-- ARR_TIME: long (nullable = true)
|-- CANCELLED: long (nullable = true)
|-- CRS_ARR_TIME: long (nullable = true)
|-- CRS_DEP_TIME: long (nullable = true)
|-- DEP_TIME: long (nullable = true)
|-- DEST: string (nullable = true)
|-- DEST_CITY_NAME: string (nullable = true)
|-- DISTANCE: long (nullable = true)
|-- FL_DATE: string (nullable = true)
|-- OP_CARRIER: string (nullable = true)
|-- OP_CARRIER_FL_NUM: long (nullable = true)
|-- ORIGIN: string (nullable = true)
|-- ORIGIN_CITY_NAME: string (nullable = true)
|-- TAXI_IN: long (nullable = true)
|-- WHEELS_ON: long (nullable = true)


```

Command took 0.13 seconds -- by prashant@scholarnest.com at 4/26/2022, 11:34:20 PM on demo-cluster

Schema inference is not reliable.

We have the following problems in schema inference:

1. Automatic inferring of schema is often incorrect
2. Inferring schema is additional work for Spark, and it takes some extra time
3. Schema inference is conflicting with the schema validation
4. It might also change the column order

Look at the schema of the data frame we created. You can see that Spark automatically inferred the schema for this Dataframe. Some columns are long, and others are string. And that is not correct.

07-dataframe-schema Python

demo-cluster (1) Spark Jobs

Command took 20.42 seconds -- by prashant@scholarnest.com at 4/26/2022, 11:22:50 PM on demo-cluster

Cmd 3

```
> 1 flight_time_raw_df.printSchema()

root
|-- ARR_TIME: long (nullable = true)
|-- CANCELLED: long (nullable = true)
|-- CRS_ARR_TIME: long (nullable = true)
|-- CRS_DEP_TIME: long (nullable = true)
|-- DEP_TIME: long (nullable = true)
|-- DEST: string (nullable = true)
|-- DEST_CITY_NAME: string (nullable = true)
|-- DISTANCE: long (nullable = true)
|-- FL_DATE: string (nullable = true)
|-- OP_CARRIER: string (nullable = true)
|-- OP_CARRIER_FL_NUM: long (nullable = true)
|-- ORIGIN: string (nullable = true)
|-- ORIGIN_CITY_NAME: string (nullable = true)
|-- TAXI_IN: long (nullable = true)
|-- WHEELS_ON: long (nullable = true)


```

Command took 0.13 seconds -- by prashant@scholarnest.com at 4/26/2022, 11:34:20 PM on demo-cluster

Schema inference may also change the order of the columns. You can see it here if you compare the schema with your data file. The first column in Dataframe is ARR\_TIME, whereas the first column in the file is FL\_DATE. Column reordering is not a big problem. We should avoid creating a logic that depends on the column ordering but changing order often causes confusion. So we must avoid it.

The screenshot shows a Jupyter Notebook interface with a Python kernel. The title bar says "07-dataframe-schema".

The code cell contains:

```
1 flight_time_raw_df.printSchema()
```

The output shows the schema:

```
root
+--- ARR_TIME: long (nullable = true)
|   |-- CANCELLED: long (nullable = true)
|   |-- CRS_ARR_TIME: long (nullable = true)
|   |-- CRS_DEP_TIME: long (nullable = true)
|   |-- DEP_TIME: long (nullable = true)
|   |-- DEST: string (nullable = true)
|   |-- DEST_CITY_NAME: string (nullable = true)
|   |-- DISTANCE: long (nullable = true)
|   |-- FL_DATE: string (nullable = true)
|   |-- OP_CARRIER: string (nullable = true)
|   |-- OP_CARRIER_FL_NUM: long (nullable = true)
|   |-- ORIGIN: string (nullable = true)
|   |-- ORIGIN_CITY_NAME: string (nullable = true)
|   |-- TAXI_IN: long (nullable = true)
|   |-- WHEELS_ON: long (nullable = true)
```

A red box highlights the first few rows of the data frame. An arrow points from the schema's ARR\_TIME field to the data's FL\_DATE field. The data rows are:

```
{
  "FL_DATE": "1/1/2000",
  "OP_CARRIER": "DL",
  "OP_CARRIER_FL_NUM": 1451,
  "ORIGIN": "BOS",
  "ORIGIN_CITY_NAME": "Boston, MA",
  "DEST": "ATL",
  "DEST_CITY_NAME": "Atlanta, GA",
  "CRS_DEP_TIME": 1115,
  "DEP_TIME": 1113,
  "WHEELS_ON": 1343,
  "TAXI_IN": 5,
  "CRS_ARR_TIME": 1400,
  "ARR_TIME": 1348,
  "CANCELLED": 0,
  "DISTANCE": 946}
```

We have two steps to enforce a schema:

1. Create a schema definition for your data file.
2. Explicitly enforce the schema

The first step is to create a schema definition.

We have two approaches to do it.

1. Schema DDL String
2. StructType Object

Schema DDL string is the easiest method.

Here is my schema definition for the flight-time data. A schema DDL is a comma-separated string of column names and data types. So the first column name is FL\_DATE, and it is a date. Similarly, the second column name is OP\_CARRIER, and it is a string.

You can create a list of all the columns and their data types.

07-dataframe-schema Python

demo-cluster Cmd 1

```
1 %fs ls /FileStore/tables
```

	path	name	size	modificationTime
1	dbfs:/FileStore/tables/flight_time.json	flight_time.json	87150442	1650995258000

Showing all 1 rows.

Command took 5.95 seconds -- by prashant@scholarnest.com at 4/26/2022, 11:20:10 PM on demo-cluster

Cmd 2

```
1 flight_schema_ddl = """FL_DATE DATE, OP_CARRIER STRING, OP_CARRIER_FL_NUM INT, ORIGIN STRING,  
2 ORIGIN_CITY_NAME STRING, DEST STRING, DEST_CITY_NAME STRING, CRS_DEP_TIME INT, DEP_TIME INT,  
3 WHEELS_ON INT, TAXI_IN INT, CRS_ARR_TIME INT, ARR_TIME INT, CANCELLED BOOLEAN, DISTANCE INT""""
```

Once a schema is defined, you can enforce it to the DataframeReader by adding a schema method to your DataframeReader.

So here I am telling the DataframeReader to read a JSON format file using the flight\_schema\_ddl.

The screenshot shows a Jupyter Notebook cell titled "07-dataframe-schema" in Python mode. The cell contains the following code:

```
1 flight_schema_ddl = """FL_DATE DATE, OP_CARRIER STRING, OP_CARRIER_FL_NUM INT, ORIGIN STRING,  
2 ORIGIN_CITY_NAME STRING, DEST STRING, DEST_CITY_NAME STRING, CRS_DEP_TIME INT, DEP_TIME INT,  
3 WHEELS_ON INT, TAXI_IN INT, CRS_ARR_TIME INT, ARR_TIME INT, CANCELLED BOOLEAN, DISTANCE INT"""
```

Below the code, a message indicates the command took 0.05 seconds. The cell is labeled "Cmd 3".

The next cell, labeled "Cmd 4", contains the following code:

```
1 flight_time_raw_df = spark.read \  
2 .format("json") \  
3 .schema(flight_schema_ddl) \  
4 .load("/FileStore/tables/flight_time.json")
```

Below this code, a message indicates the command took 25.65 seconds. The cell is labeled "Cmd 4".

You can re-check the schema.  
The print schema is now showing it correctly.

07-dataframe-schema Python

demo-cluster

Cmd 4

```
1 flight_time_raw_df.printSchema()
```

root

```
|-- FL_DATE: date (nullable = true)
|-- OP_CARRIER: string (nullable = true)
|-- OP_CARRIER_FL_NUM: integer (nullable = true)
|-- ORIGIN: string (nullable = true)
|-- ORIGIN_CITY_NAME: string (nullable = true)
|-- DEST: string (nullable = true)
|-- DEST_CITY_NAME: string (nullable = true)
|-- CRS_DEP_TIME: integer (nullable = true)
|-- DEP_TIME: integer (nullable = true)
|-- WHEELS_ON: integer (nullable = true)
|-- TAXI_IN: integer (nullable = true)
|-- CRS_ARR_TIME: integer (nullable = true)
|-- ARR_TIME: integer (nullable = true)
|-- CANCELLED: boolean (nullable = true)
|-- DISTANCE: integer (nullable = true)
```

Command took 0.12 seconds -- by prashant@scholarnest.com at 4/27/2022, 5:36:41 PM on demo-cluster

We can check the data once and see if we loaded it correctly.

So here is my Dataframe shown below. It looks like FL\_DATE is not read correctly, and hence it is NULL. I do see some nulls in between, but other rows have values.

07-dataframe-schema Python

```
demo-cluster
|-- ARR_DELAY: integer (nullable = true)
|-- CANCELLED: boolean (nullable = true)
|-- DISTANCE: integer (nullable = true)

Command took 0.12 seconds -- by prashant@scholarnest.com at 4/27/2022, 5:36:41 PM on demo-cluster
```

Cmd 5

```
1 display(flight_time_raw_df)
```

▶ (1) Spark Jobs

	FL_DATE	OP_CARRIER	OP_CARRIER_FL_NUM	ORIGIN	ORIGIN_CITY_NAME	DEST	DEST_CITY_NAME	CRS_DEP_TIME	DEP_TIME
1	null	DL	1451	BOS	Boston, MA	ATL	Atlanta, GA	1115	1113
2	null	DL	1479	BOS	Boston, MA	ATL	Atlanta, GA	1315	1311
3	null	DL	1857	BOS	Boston, MA	ATL	Atlanta, GA	1415	1414
4	null	DL	1997	BOS	Boston, MA	ATL	Atlanta, GA	1715	1720
5	null	DL	2065	BOS	Boston, MA	ATL	Atlanta, GA	2015	2010
6	null	US	2619	BOS	Boston, MA	ATL	Atlanta, GA	650	649
7	null	US	2621	BOS	Boston, MA	ATL	Atlanta, GA	1440	1446

Truncated results, showing first 1000 rows.  
Click to re-execute with maximum result limits.

We can check the data once and see if we loaded it correctly.  
So here is my Dataframe shown below. It looks like FL\_DATE is not read correctly, and hence it is NULL. I do see some nulls in between, but other rows have values. The CANCELLED column is also null everywhere.

07-dataframe-schema Python

demo-cluster

|-- ARR\_TIME: integer (nullable = true)  
|-- CANCELLED: boolean (nullable = true)  
|-- DISTANCE: integer (nullable = true)

Command took 0.12 seconds -- by prashant@scholarnest.com at 4/27/2022, 5:36:41 PM on demo-cluster

Cmd 5

```
1 display(flight_time_raw_df)
```

(1) Spark Jobs

FL_DATE	OP_CARRIER	OP_CARRIER_FL_NUM	ORIGIN	ORIGIN_CITY_NAME	DEST	DEST_CITY_NAME	CRS_DEP_TIME	DEP_TIME
null	DL	1451	BOS	Boston, MA	ATL	Atlanta, GA	1115	1113
null	DL	1479	BOS	Boston, MA	ATL	Atlanta, GA	1315	1311
null	DL	1857	BOS	Boston, MA	ATL	Atlanta, GA	1415	1414
null	DL	1997	BOS	Boston, MA	ATL	Atlanta, GA	1715	1720
null	DL	2065	BOS	Boston, MA	ATL	Atlanta, GA	2015	2010
null	US	2619	BOS	Boston, MA	ATL	Atlanta, GA	650	649
null	US	2621	BOS	Boston, MA	ATL	Atlanta, GA	1440	1446

Truncated results, showing first 1000 rows.  
Click to re-execute with maximum result limits.

Here is my schema definition shown below.

The FL\_DATE is a date, and the CANCELLED column is a BOOLEAN.

The date, timestamp, and the boolean values are often loaded incorrectly, or they became null. You will see it almost everywhere.

### 07-dataframe-schema

Python

demo-cluster

Cmd 4

```
1 flight_time_raw_df.printSchema()

root
|-- FL_DATE: date (nullable = true)
|-- OP_CARRIER: string (nullable = true)
|-- OP_CARRIER_FL_NUM: integer (nullable = true)
|-- ORIGIN: string (nullable = true)
|-- ORIGIN_CITY_NAME: string (nullable = true)
|-- DEST: string (nullable = true)
|-- DEST_CITY_NAME: string (nullable = true)
|-- CRS_DEP_TIME: integer (nullable = true)
|-- DEP_TIME: integer (nullable = true)
|-- WHEELS_ON: integer (nullable = true)
|-- TAXI_IN: integer (nullable = true)
|-- CRS_ARR_TIME: integer (nullable = true)
|-- ARR_TIME: integer (nullable = true)
|-- CANCELLED: boolean (nullable = true)
|-- DISTANCE: integer (nullable = true)
```

Command took 0.12 seconds -- by prashant@scholarnest.com at 4/27/2022, 5:36:41 PM on demo-cluster

If you check the requirements, you can see that FL\_DATE is the flight date in MM/dd/yyyy format, and the CANCELLED is one or zero.

The default data format for DataFramereader is yyyy-MM-dd.

So the DataFramereader expected the FL\_DATE to be in yyyy-MM-dd.

Similarly, the default format for a BOOLEAN is TRUE or FALSE.

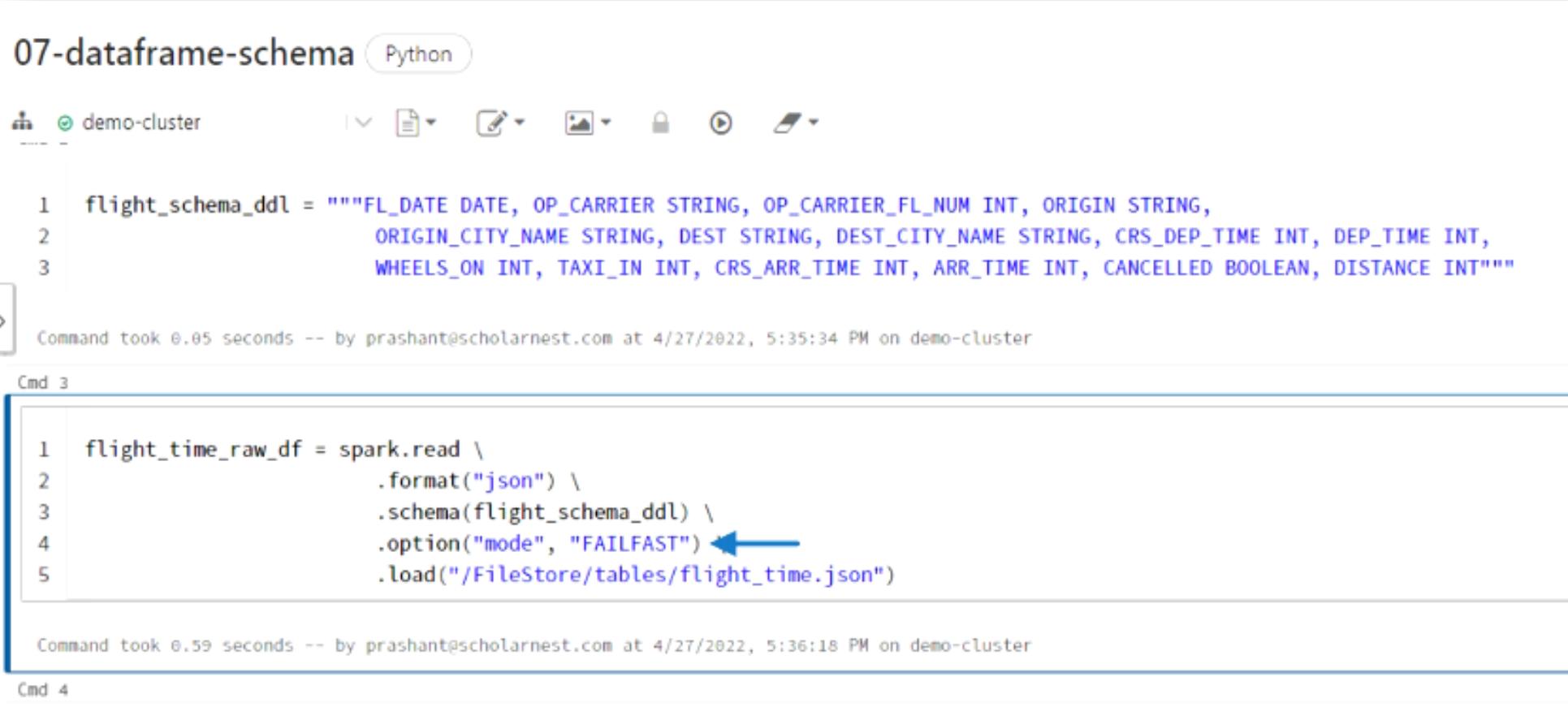
So the DataFrameReader expected the FL\_DATE to be TRUE or FALSE.

When the data is not found in the expected format, it fails the schema validation, and the Spark reads a null.

It doesn't throw an error. However, you can change this behaviour by setting the mode option.

So am setting the DataFrameReader mode to FAILFAST.

So now, if the DataFrameReader detects an incorrect value, it will stop and throw an exception.



The screenshot shows a Jupyter Notebook cell titled "07-dataframe-schema" in Python mode. The cell contains the following code:

```
1 flight_schema_ddl = """FL_DATE DATE, OP_CARRIER STRING, OP_CARRIER_FL_NUM INT, ORIGIN STRING,  
2 ORIGIN_CITY_NAME STRING, DEST STRING, DEST_CITY_NAME STRING, CRS_DEP_TIME INT, DEP_TIME INT,  
3 WHEELS_ON INT, TAXI_IN INT, CRS_ARR_TIME INT, ARR_TIME INT, CANCELLED BOOLEAN, DISTANCE INT"""
```

Below the code, a message indicates the command took 0.05 seconds. The cell is labeled "Cmd 3".

Below the cell, another command is shown:

```
1 flight_time_raw_df = spark.read \  
2 .format("json") \  
3 .schema(flight_schema_ddl) \  
4 .option("mode", "FAILFAST") ←  
5 .load("/FileStore/tables/flight_time.json")
```

This command also took 0.59 seconds. The cell is labeled "Cmd 4". A blue arrow points to the ".option('mode', 'FAILFAST')" line.

The DataFrameReader mode option can take three values.

1. FAIL FAST

The fail-fast will throw an error when detecting a malformed record.

2. DROPMALFORMED

The DROPMALFORMED will leave the malformed record and load others.

3. PERMISSIVE

The PERMISSIVE is the default value, and it will set the malformed column to null.

If you try running the display. You see an exception as shown below.

But how do we fix it?

Well, the error is coming because the data is not in the expected format.

We cannot change the file data. But we can change the expected format.

The screenshot shows a Jupyter Notebook interface with the title "07-dataframe-schema". A Python code cell (Cmd 5) contains the command `display(flight_time_raw_df)`. The output of this command is an error message:

```
FileReadException: Error while reading file dbfs:/FileStore/tables/flight_time.json.  
Caused by: SparkException: Malformed records are detected in record parsing. Parse Mode: FAILFAST. To process malformed records as null result, try setting the option 'mode' as 'PERMISSIVE'.  
Caused by: BadRecordException: java.time.DateTimeException: Cannot cast 1/1/2000 to DateType. To return NULL instead, use 'try_cast'. If necessary set spark.sql.ansi.enabled to false to bypass this error.  
Caused by: DateTimeException: Cannot cast 1/1/2000 to DateType. To return NULL instead, use 'try_cast'. If necessary set spark.sql.ansi.enabled to false to bypass this error.
```

The error message is preceded by "Command took 0.12 seconds -- by prashant@scholarnest.com at 4/27/2022, 5:36:41 PM on demo-cluster".

A red vertical bar highlights the error message, and a tooltip "Error running command 5. Go to command" appears over the bar. The status bar at the bottom shows "Cmd 6".

I am adding another DataFrameReader option.

I am setting the expected date format to M/d/y because my file data is in the M/d/y format.

With this change, the FL\_DATE column will load correctly.

Cmd 3

```
1 flight_time_raw_df = spark.read \
2     .format("json") \
3     .schema(flight_schema_ddl) \
4     .option("mode", "FAILFAST") \
5     .option("dateFormat", "M/d/y") ←
6     .load("/FileStore/tables/flight_time.json")
```

Command took 0.59 seconds -- by prashant@scholarnest.com at 4/27/2022, 5:47:00 PM on demo-cluster

Unfortunately, the DataFrameReader does not allow us to change the expected Boolean format. So I do not have an option to load the zero/one as Boolean.

But that's not a problem, we have a straightforward approach. Set the columns to string and load it. We can convert it to appropriate data types after loading.

So I have changed the CANCELLED column to a string so we can load it.

Cmd 2

```
1 flight_schema_ddl = """FL_DATE DATE, OP_CARRIER STRING, OP_CARRIER_FL_NUM INT, ORIGIN STRING,  
2 ORIGIN_CITY_NAME STRING, DEST STRING, DEST_CITY_NAME STRING, CRS_DEP_TIME INT, DEP_TIME INT,  
3 WHEELS_ON INT, TAXI_IN INT, CRS_ARR_TIME INT, ARR_TIME INT, CANCELLED STRING, DISTANCE INT"""
```

Command took 0.06 seconds -- by prashant@scholarnest.com at 4/27/2022, 5:47:40 PM on demo-cluster

Cmd 3

```
1 flight_time_raw_df = spark.read \  
2 .format("json") \  
3 .schema(flight_schema_ddl) \  
4 .option("mode", "FAILFAST") \  
5 .option("dateFormat", "M/d/y") \  
6 .load("/FileStore/tables/flight_time.json")
```

Command took 0.59 seconds -- by prashant@scholarnest.com at 4/27/2022, 5:47:00 PM on demo-cluster

Python

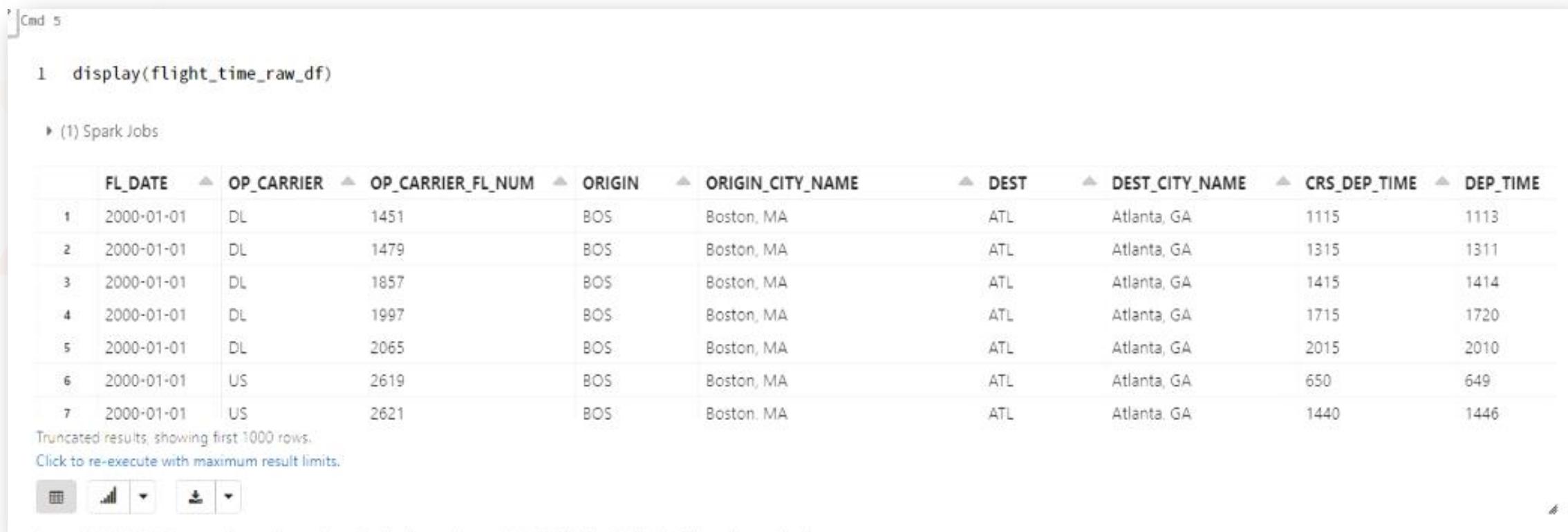
Now if you try to check the Dataframe, you can see that it is loaded correctly as shown below.

```
Cmd 5
1 display(flight_time_raw_df)

▶ (1) Spark Jobs
```

	FL_DATE	OP_CARRIER	OP_CARRIER_FL_NUM	ORIGIN	ORIGIN_CITY_NAME	DEST	DEST_CITY_NAME	CRS_DEP_TIME	DEP_TIME
1	2000-01-01	DL	1451	BOS	Boston, MA	ATL	Atlanta, GA	1115	1113
2	2000-01-01	DL	1479	BOS	Boston, MA	ATL	Atlanta, GA	1315	1311
3	2000-01-01	DL	1857	BOS	Boston, MA	ATL	Atlanta, GA	1415	1414
4	2000-01-01	DL	1997	BOS	Boston, MA	ATL	Atlanta, GA	1715	1720
5	2000-01-01	DL	2065	BOS	Boston, MA	ATL	Atlanta, GA	2015	2010
6	2000-01-01	US	2619	BOS	Boston, MA	ATL	Atlanta, GA	650	649
7	2000-01-01	US	2621	BOS	Boston, MA	ATL	Atlanta, GA	1440	1446

Truncated results, showing first 1000 rows.  
Click to re-execute with maximum result limits.



We have two approaches to defining a Schema.

1. Schema DDL String
2. StructType Object

We saw the first approach.

Now let's define the schema using the StructType object definition.

We have to import Spark types from the `pyspark.sql.types` package.

Then I can define a variable and set it to `StructType()`. The `StructType()` method takes a list of `StructField()`. You can put a comma and add all the available fields as shown below.

07-dataframe-schema Python

demo-cluster Cmd: 6

```
1 from pyspark.sql.types import *
2
3 flight_schema = StructType([
4     StructField("FL_DATE", DateType()),
5     StructField("OP_CARRIER", StringType()),
6     StructField("OP_CARRIER_FL_NUM", IntegerType()),
7     StructField("ORIGIN", StringType()),
8     StructField("ORIGIN_CITY_NAME", StringType()),
9     StructField("DEST", StringType()),
10    StructField("DEST_CITY_NAME", StringType()),
11    StructField("CRS_DEP_TIME", IntegerType()),
12    StructField("DEP_TIME", IntegerType()),
13    StructField("WHEELS_ON", IntegerType()),
14    StructField("TAXI_IN", IntegerType()),
15    StructField("CRS_ARR_TIME", IntegerType()),
16    StructField("ARR_TIME", IntegerType()),
17    StructField("CANCELLED", IntegerType()),
18    StructField("DISTANCE", IntegerType())
19])
```

Command took 0.07 seconds -- by prashant@scholarnest.com at 4/27/2022, 7:21:55 PM on demo-cluster

We use the StructType object in the same way as we used the DDL schema. Here is the code for the same. So, you can define the schema using a DDL string or using StructType.

Cmd 7

```
1 flight_time_raw_df_1 = spark.read \
2         .format("json") \
3         .schema(flight_schema) \
4         .option("mode", "FAILFAST") \
5         .option("dateFormat", "M/d/y") \
6         .load("/FileStore/tables/flight_time.json")
```

Command took 0.51 seconds -- by prashant@scholarnest.com at 4/27/2022, 7:23:59 PM on demo-cluster



Thank You  
ScholarNest Technologies Pvt Ltd.  
[www.scholarnest.com](http://www.scholarnest.com)



ScholarNest

# Spark Azure Databricks

Databricks Spark Certification and beyond

Instructor: Prashant Kumar Pandey



Absolute Beginner to Specialization in Apache Spark and Azure Databricks





# Writing Spark Dataframes

# Requirement:

In the previous lecture, we saw the following requirement.

We have a JSON data set, and you are asked to load this dataframe and create a Spark table. For your reference, here is one sample record, and here are details about the columns shown below.

## Flight Delay Data

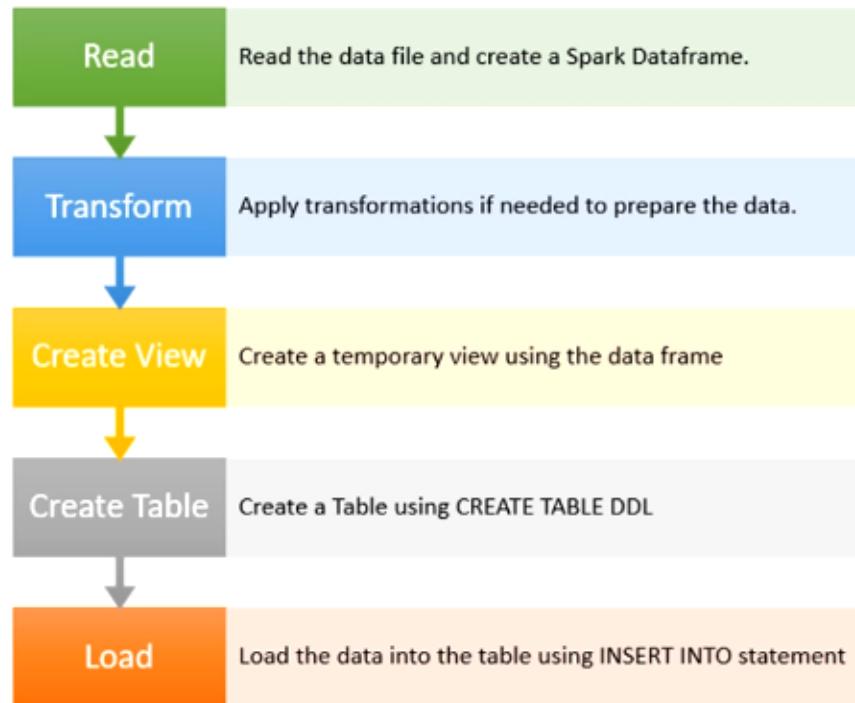
```
{  
    "FL_DATE": "1/1/2000",  
    "OP_CARRIER": "DL",  
    "OP_CARRIER_FL_NUM": 1451,  
    "ORIGIN": "BOS",  
    "ORIGIN_CITY_NAME": "Boston, MA",  
    "DEST": "ATL",  
    "DEST_CITY_NAME": "Atlanta, GA",  
    "CRS_DEP_TIME": 1115,  
    "DEP_TIME": 1113,  
    "WHEELS_ON": 1343,  
    "TAXI_IN": 5,  
    "CRS_ARR_TIME": 1400,  
    "ARR_TIME": 1348,  
    "CANCELLED": 0,  
    "DISTANCE": 946  
}
```

1. Flight date in MM/DD/YYYY format
2. Unique carrier code
3. Flight number
4. Origin IATA Airport Code
5. Origin City Name
6. Destination IATA Airport Code
7. Destination City Name
8. Scheduled Departure time in HHMM
9. Actual Departure time in HHMM
10. Flight Duration in HHMM
11. TaxiIn taxi in time, in minutes
12. Scheduled Arrival time in HHMM
13. Actual Departure time in HHMM
14. Was it canceled? 1 for a true, and 0 is false
15. Distance in miles

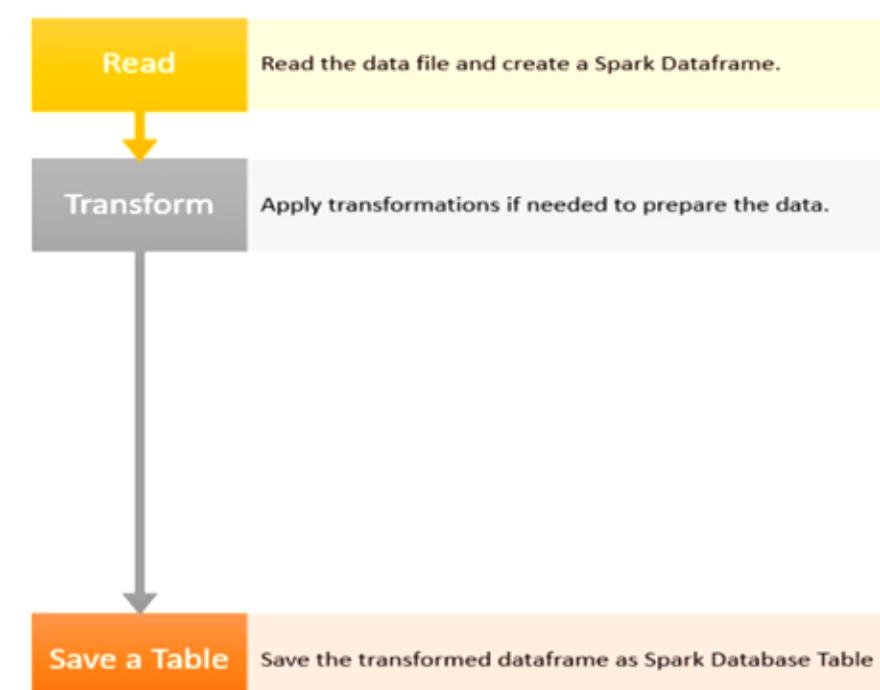
In order to serve our requirement, we have the following two approaches.

We saw the first approach in the previous documentation, now let us see the second approach in this documentation.

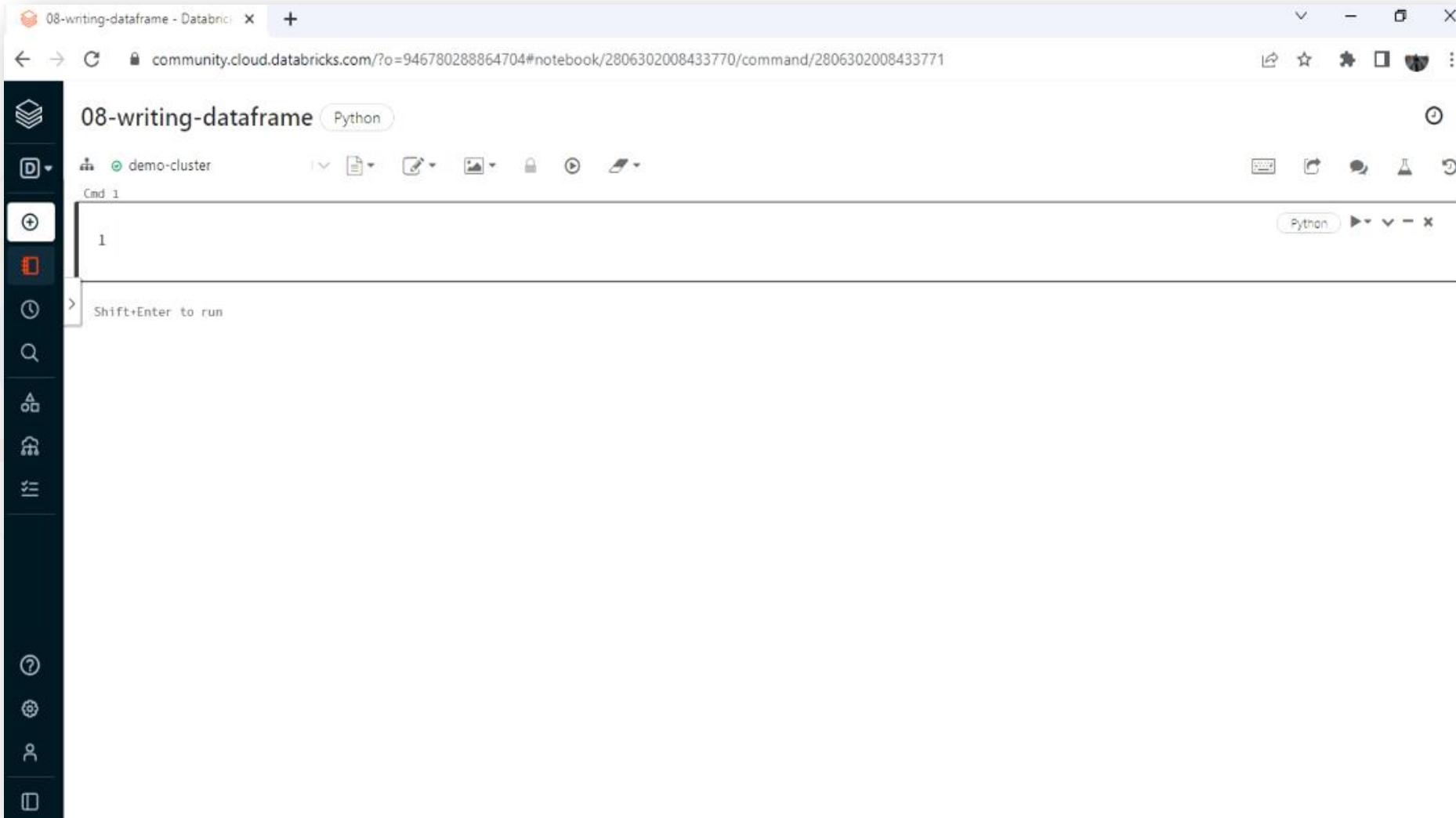
### Approach - 1



### Approach - 2



Go to your Databricks Workspace and create a new notebook. (**Reference - 08-writing-dataframe.ipynb**)



We already wrote the code for the following:

1. Define a flight delay dataset schema.
2. Load data from the file and create a Dataframe.

So I have copy-pasted the code here.

The screenshot shows a Jupyter Notebook interface with the title "08-writing-dataframe" and a Python tab selected. The interface includes a toolbar with icons for file operations, a cluster status indicator ("demo-cluster"), and a cell navigation bar. There are two code cells:

**Cmd 1**

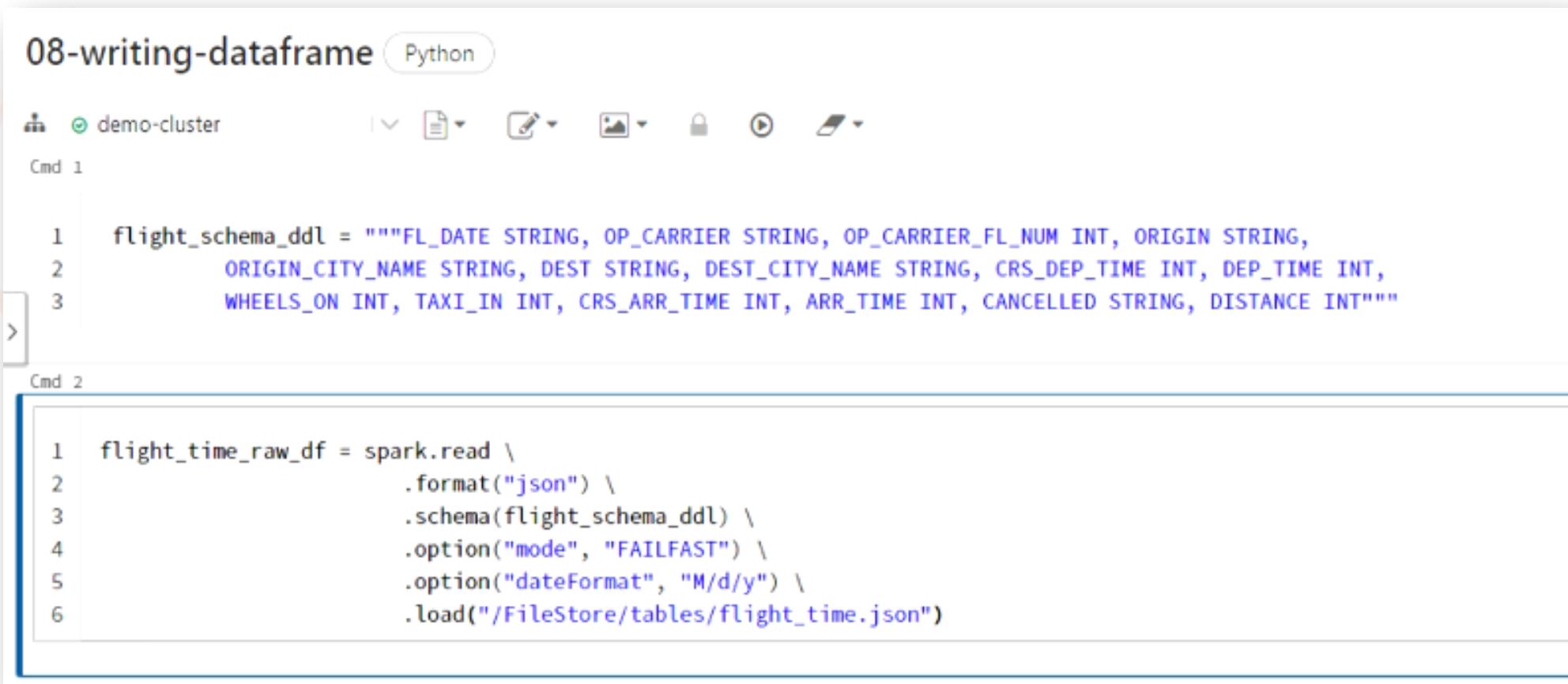
```
1 flight_schema_ddl = """FL_DATE STRING, OP_CARRIER STRING, OP_CARRIER_FL_NUM INT, ORIGIN STRING,  
2     ORIGIN_CITY_NAME STRING, DEST STRING, DEST_CITY_NAME STRING, CRS_DEP_TIME INT, DEP_TIME INT,  
3     WHEELS_ON INT, TAXI_IN INT, CRS_ARR_TIME INT, ARR_TIME INT, CANCELLED STRING, DISTANCE INT"""
```

**Cmd 2**

```
1 flight_time_raw_df = spark.read \  
2     .format("json") \  
3     .schema(flight_schema_ddl) \  
4     .option("mode", "FAILFAST") \  
5     .option("dateFormat", "M/d/y") \  
6     .load("/FileStore/tables/flight_time.json")
```

So we have this `flight_time_raw_df`. And I want to save this dataframe into a Spark table. That is my requirement.

But before saving this data as a Spark table, I want to correct the `FL_DATE` and `CANCELLED` data types.



The screenshot shows a Jupyter Notebook interface with the title "08-writing-dataframe" and the Python tab selected. The interface includes a toolbar with icons for file operations, a cluster status indicator ("demo-cluster"), and a command history area.

**Cmd 1:**

```
1 flight_schema_ddl = """FL_DATE STRING, OP_CARRIER STRING, OP_CARRIER_FL_NUM INT, ORIGIN STRING,  
2     ORIGIN_CITY_NAME STRING, DEST STRING, DEST_CITY_NAME STRING, CRS_DEP_TIME INT, DEP_TIME INT,  
3     WHEELS_ON INT, TAXI_IN INT, CRS_ARR_TIME INT, ARR_TIME INT, CANCELLED STRING, DISTANCE INT"""
```

**Cmd 2:**

```
1 flight_time_raw_df = spark.read \  
2     .format("json") \  
3     .schema(flight_schema_ddl) \  
4     .option("mode", "FAILFAST") \  
5     .option("dateFormat", "M/d/y") \  
6     .load("/FileStore/tables/flight_time.json")
```

We have already learned about the following:

1. `withColumn()` transformations
2. `to_date()` function

In fact, we converted a string field to a date field in an earlier lecture. So I guess you can fix the FL\_DATE quickly.

But how do you fix the CANCELLED?

The idea is to look for something that helps you write the following logic:

`if CANCELLED == 1 then true else false`

Functions in Spark are defined in two places.

1. `pyspark.sql.functions` package
2. Built-In Spark SQL functions

You can check for the Built-In SQL functions in the documentation, and you will see an `if` function, and it takes three arguments.

The first argument is an expression that must evaluate as true or false.

The other two arguments are true or false.

I can use this function in my Dataframe code to transform the CANCELLED column from a number to a boolean.

Here is the code to transform the FL\_DATE and CANCELLED columns.

I am importing to\_date() and expr() functions because I need them. Then I am starting with flight\_time\_raw\_df and transforming two columns using the withColumn() transformation.

The withColumn() takes two arguments. The first argument is the column's name that you want to transform. The second argument is the transformation logic. The logic for FL\_DATE is as simple as applying the to\_date() function, so we change it from a string to a date.

The logic for the CANCELLED column is also straightforward. We want to apply the if() built-in function as shown here. If CANCELLED is one, make it True else, make it False. Also, my logic for transforming the CANCELLED column is a SQL expression. It is not a dataframe expression. So I must explicitly evaluate the SQL expression. And that's why we use the expr() function.

Cmd 3

```
1 from pyspark.sql.functions import to_date, expr  
2  
3 flight_time_df = flight_time_raw_df \  
4         .withColumn("FL_DATE", to_date("FL_DATE", "M/d/y")) \  
5         .withColumn("CANCELLED", expr("if(CANCELLED=1, true, false)"))
```

Command took 0.09 seconds -- by prashant@scholarnest.com at 4/29/2022, 1:05:17 AM on demo-cluster

Cmd 4

You can save your Dataframe into a database table using the code shown below.

Cmd 4

```
1 flight_time_df.write \
2         .format("parquet") \
3         .mode("overwrite") \
4         .saveAsTable("flight_time_tbl")|
```

If you look at the code, it looks similar to the DataFrameReader. So you can see the spark.read here. And you can see the dataframe.write here. The first code block is reading data from a source and creating a Dataframe. The second code block is from writing Dataframe to a target. The first one starts with a spark session and access the DataFrameReader using the read attribute. The second one starts with a Dataframe and access the DataFrameWriter using the write attribute. Once you have a DataFrameWriter, the rest of the code structure is the same as the DataFrameReader code. We can set the format, set some options, and save it.

08-writing-dataframe Python

demo-cluster Cmd 2

```
1 flight_time_raw_df = spark.read \
2           .format("json") \
3           .schema(flight_schema_ddl) \
4           .option("mode", "FAILFAST") \
5           .option("dateFormat", "M/d/y") \
6           .load("/FileStore/tables/flight_time.json")
```

Command took 0.49 seconds -- by prashant@scholarnest.com at 4/29/2022, 1:00:28 AM on demo-cluster

Cmd 3

```
1 from pyspark.sql.functions import to_date, expr
2
3 flight_time_df = flight_time_raw_df \
4           .withColumn("FL_DATE", to_date("FL_DATE", "M/d/y")) \
5           .withColumn("CANCELLED", expr("if(CANCELLED==1, true, false)"))
```

Command took 0.09 seconds -- by prashant@scholarnest.com at 4/29/2022, 1:05:17 AM on demo-cluster

Cmd 4

```
1 flight_time_df.write \
2           .format("parquet") \
3           .mode("overwrite") \
4           .saveAsTable("flight_time_tbl")
```

DataFrameWriter supports the following modes:

1. append: Append contents of this DataFrame to existing data.
2. overwrite: Overwrite existing data.
3. error or errorIfExists: Throw an exception if data already exists.
4. ignore: Silently ignore this operation if data already exists.

If you look at the data section, you will see a table is created for your Dataframe.

The screenshot shows the Databricks Data Catalog interface. On the left, the navigation sidebar is visible with options like 'Data Science & En...', 'Create', 'Workspace', 'Recents', 'Search', 'Data' (which is selected and highlighted in blue), 'Compute', 'Jobs', 'Help', 'Settings', and 'Menu options'. The main area is divided into two tabs: 'Database Tables' (selected) and 'DBFS'. Under 'Database Tables', there is a 'Databases' section with a dropdown menu set to 'default', and a 'Tables' section. A table named 'flight\_time\_tbl' is listed in the 'Tables' section. To the right of the table name, a portion of the table schema is shown in JSON format:  `"false"))`. At the top right of the main area, there is a 'Create Table' button and a small icon.

You can also describe the table from your notebook as shown below. And you will see that both FL\_DATE and CANCELLED columns are rectified as per your requirement.

```
1 %sql
2 describe extended flight_time_tbl
```

	col_name	data_type	comment
1	FL_DATE	date	null
2	OP_CARRIER	string	null
3	OP_CARRIER_FL_NUM	int	null
4	ORIGIN	string	null
5	ORIGIN_CITY_NAME	string	null
6	DEST	string	null
7	DEST CITY NAME	string	null

Showing all 30 rows.

Command took 0.48 seconds -- by prashant@scholarnest.com at 4/29/2022, 1:09:48 AM on demo-cluster



Thank You  
ScholarNest Technologies Pvt Ltd.  
[www.scholarnest.com](http://www.scholarnest.com)



ScholarNest

# Spark Azure Databricks

Databricks Spark Certification and beyond

Instructor: Prashant Kumar Pandey



Absolute Beginner to Specialization in Apache Spark and Azure Databricks





# Managed Vs External Tables

**Spark allows you to create two types of tables.**

- 1. Managed Tables**
- 2. External Tables**

Go to your databricks workspace and open an existing workbook (**Reference - 08-writing-dataframe**).

In this example, we learned to create a database table.

The screenshot shows a Databricks notebook interface with the title "08-writing-dataframe". The notebook contains three code cells:

- Cmd 1:** Contains Python code to define a flight schema DDL:

```
flight_schema_ddl = """FL_DATE STRING, OP_CARRIER STRING, OP_CARRIER_FL_NUM INT, ORIGIN STRING,  
ORIGIN_CITY_NAME STRING, DEST STRING, DEST_CITY_NAME STRING, CRS_DEP_TIME INT, DEP_TIME INT,  
WHEELS_ON INT, TAXI_IN INT, CRS_ARR_TIME INT, ARR_TIME INT, CANCELLED STRING, DISTANCE INT"""
```

Command took 0.04 seconds -- by prashant@scholarnest.com at 4/29/2022, 1:00:25 AM on demo-cluster
- Cmd 2:** Contains Python code to read JSON data into a DataFrame:

```
flight_time_raw_df = spark.read \  
.format("json") \  
.schema(flight_schema_ddl) \  
.option("mode", "FAILFAST") \  
.option("dateFormat", "M/d/y") \  
.load("/FileStore/tables/flight_time.json")
```

Command took 0.49 seconds -- by prashant@scholarnest.com at 4/29/2022, 1:00:28 AM on demo-cluster
- Cmd 3:** Contains Python code to transform the raw DataFrame:

```
from pyspark.sql.functions import to_date, expr  
  
flight_time_df = flight_time_raw_df \  
.withColumn("FL_DATE", to_date("FL_DATE", "M/d/y")) \  
.withColumn("CANCELLED", expr("if(CANCELLED==1, true, false)"))
```

Command took 0.09 seconds -- by prashant@scholarnest.com at 4/29/2022, 1:05:17 AM on demo-cluster

# Click the file menu and clone your notebook as shown below.

The screenshot shows a Jupyter Notebook interface with the following details:

- Title Bar:** 08-writing-dataframe (Python)
- Cluster Selection:** demo-cluster
- Toolbar:** Includes icons for New Notebook, Clone (highlighted), Rename, Move, Delete, Upload Data, Export, Publish, Clear Revision History, and Change Default Language.
- Cell 4 (Cmd 4):** Python code for writing a DataFrame to a Parquet file.

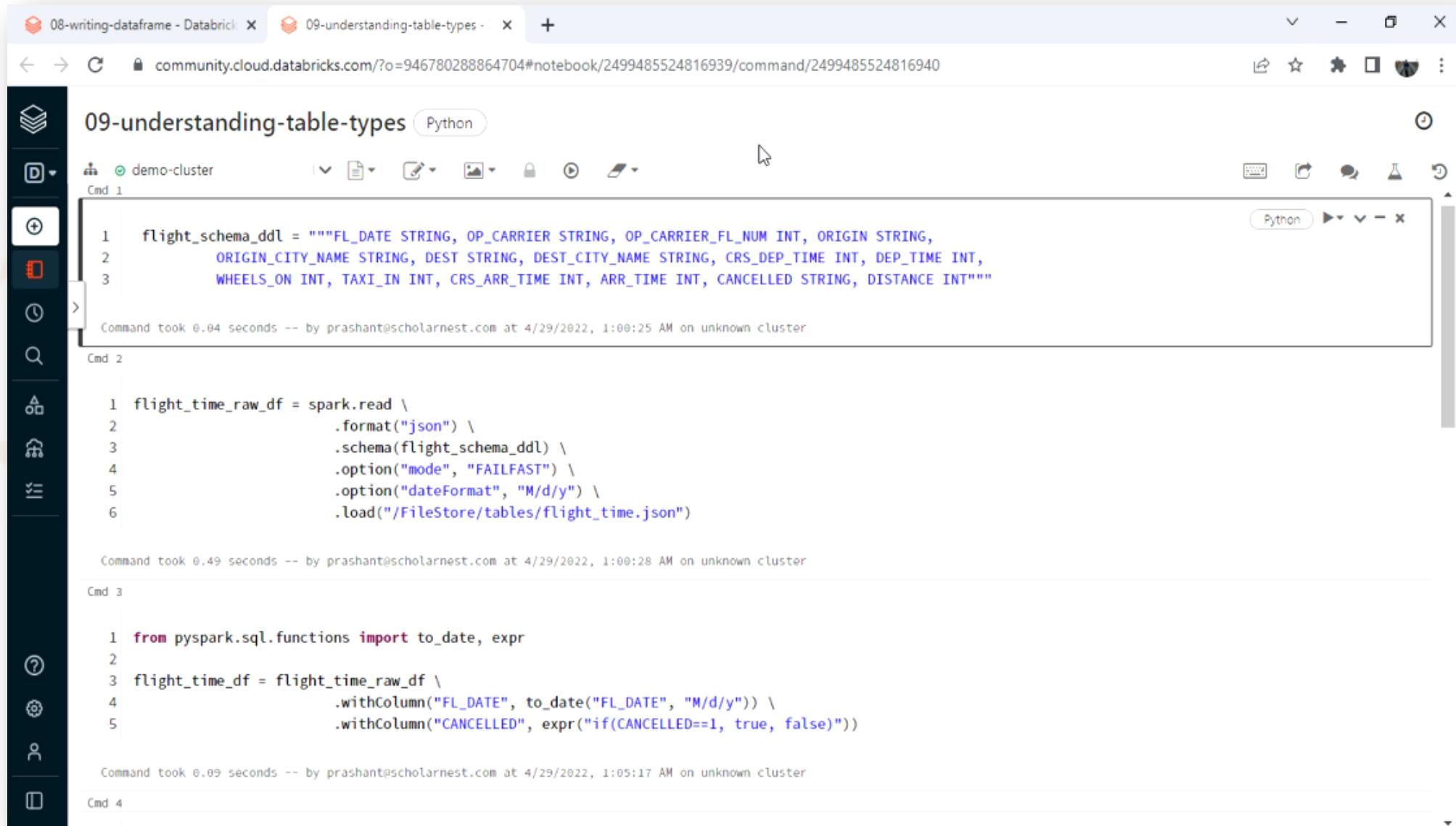
```
flight_time_df.write \  
.format("parquet") \  
.mode("overwriting") \  
.saveAsTable("flights")
```

Command took 0.09 seconds -- by p...
- Cell 5 (Cmd 5):** SQL command to describe the newly created table.

```
%sql  
describe extended flight_time_tbl
```

Command took 18.64 seconds -- by p...
- File Menu Context:** A context menu is open over the cluster selection area, listing options: New Notebook, Clone (highlighted), Rename, Move, Delete, Upload Data, Export, Publish, Clear Revision History, and Change Default Language.
- Timestamps:** 29/2022, 1:05:17 AM on demo-cluster and 29/2022, 1:07:55 AM on demo-cluster.

Once cloned, rename your notebook and we will start working on this notebook.  
**(Reference - 09-understanding-table-types.ipynb)**



The screenshot shows a Databricks notebook interface with the title "09-understanding-table-types" and the language "Python". The notebook contains the following code:

```
flight_schema_ddl = """FL_DATE STRING, OP_CARRIER STRING, OP_CARRIER_FL_NUM INT, ORIGIN STRING,  
ORIGIN_CITY_NAME STRING, DEST STRING, DEST_CITY_NAME STRING, CRS_DEP_TIME INT, DEP_TIME INT,  
WHEELS_ON INT, TAXI_IN INT, CRS_ARR_TIME INT, ARR_TIME INT, CANCELLED STRING, DISTANCE INT"""  
  
flight_time_raw_df = spark.read \  
.format("json") \  
.schema(flight_schema_ddl) \  
.option("mode", "FAILFAST") \  
.option("dateFormat", "M/d/y") \  
.load("/FileStore/tables/flight_time.json")  
  
from pyspark.sql.functions import to_date, expr  
flight_time_df = flight_time_raw_df \  
.withColumn("FL_DATE", to_date("FL_DATE", "M/d/y")) \  
.withColumn("CANCELLED", expr("if(CANCELLED==1, true, false)"))
```

Each command includes a timestamp and the user who ran it.

If you try running this notebook once. You might see an error as shown below. Your cluster was terminated, and you created a new cluster. So, you will see this error. Because the metadata and compute layers were removed by the Databricks. But the data files are still there.

The screenshot shows a Databricks notebook interface with the title '09-understanding-table-types'. The notebook is set to run in Python. The code in the notebook is as follows:

```
flight_time_df = flight_time_raw_df \
    .withColumn("FL_DATE", to_date("FL_DATE", "M/d/y")) \
    .withColumn("CANCELLED", expr("if(CANCELLED==1, true, false)"))

Command took 0.29 seconds -- by prashant@scholarnest.com at 4/29/2022, 10:16:37 PM on demo-cluster
```

Cmd 4

```
flight_time_df.write \
    .format("parquet") \
    .mode("overwrite") \
    .saveAsTable("flight_time_tbl")

AnalysisException: Can not create the managed table(``flight_time_tbl``). The associated location('dbfs:/user/hive/warehouse/flight_time_tbl') already exists.

Command took 3.71 seconds -- by prashant@scholarnest.com at 4/29/2022, 10:16:37 PM on demo-cluster
```

Cmd 5

```
%sql
describe extended flight_time_tbl

Command skipped

Command took 3.70 seconds -- by prashant@scholarnest.com at 4/29/2022, 10:16:37 PM on demo-cluster
```

Cmd 6

```
1
```

A red box highlights the error message from Cmd 4: 'AnalysisException: Can not create the managed table(``flight\_time\_tbl``). The associated location('dbfs:/user/hive/warehouse/flight\_time\_tbl') already exists.' A red vertical bar also appears near the top right of the notebook area, indicating an error.

You can fix the error using the %fs rm -r command, add a new cell above and write the %fs command as shown below.

Now, if you try to run the cells it will work perfectly fine.

09-understanding-table-types Python

demo-cluster

```
2
3 flight_time_df = flight_time_raw_df \
4         .withColumn("FL_DATE", to_date("FL_DATE", "M/d/y")) \
5         .withColumn("CANCELLED", expr("if(CANCELLED==1, true, false)"))
```

Command took 0.29 seconds -- by prashant@scholarnest.com at 4/29/2022, 10:16:37 PM on demo-cluster

Cmd 4

```
1 %fs rm -r user/hive/warehouse/
```

res0: Boolean = true

Command took 14.32 seconds -- by prashant@scholarnest.com at 4/29/2022, 10:17:47 PM on demo-cluster

Cmd 5

```
1 flight_time_df.write \
2     .format("parquet") \
3     .mode("overwrite") \
4     .saveAsTable("flight_time_tbl")
```

AnalysisException: Can not create the managed table(''flight\_time\_tbl''). The associated location('dbfs:/user/hive/warehouse/flight\_time\_tbl') already exists.

Command took 3.71 seconds -- by prashant@scholarnest.com at 4/29/2022, 10:16:37 PM on demo-cluster

Scroll down to see the output of the describe command, and you will see the table type. This table is a managed table.  
So by default, the saveAsTable() method created a managed table.  
You can also create a table using the create table DDL statement.  
That one also creates a managed table by default.



```
> | Cmd 6
1 %sql
2 describe extended flight_time_tbl
```

col_name	data_type	comment
23	Created By	Spark 3.2.1
24	Type	MANAGED
25	Provider	parquet
26	Statistics	3767904 bytes
27	Location	dbfs:/user/hive/warehouse/flight_time_tbl
28	Serde Library	org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe

Showing all 30 rows.

Command took 6.87 seconds -- by prashant@scholarnest.com at 4/29/2022, 10:18:15 PM on demo-cluster

We have two properties of a managed table:

1. Spark creates managed table at a predefined warehouse location.
2. Spark manages table metadata and table data together.

You might ask the following questions about the first point:

1. Can we define this location?
2. Can we change it?

The answer is YES-NO.

You can configure the default warehouse directory location using the `spark.SQL.warehouse.dir` configuration. You can go to Spark UI environment tab and search for the `spark.SQL.warehouse.dir` configuration.

The value for this configuration will be redacted for security reasons.

But the `spark.SQL.warehouse.dir` defines the default location for the managed tables.

You can set this configuration before starting the Spark cluster.

There are two things to learn here:

1. Spark creates all the managed tables at a fixed warehouse directory location.
2. You can set the warehouse directory location using the `spark.SQL.warehouse.dir` configuration.

We have two properties of a managed table:

1. Spark creates managed table at a predefined warehouse location.
2. Spark manages table metadata and table data together.

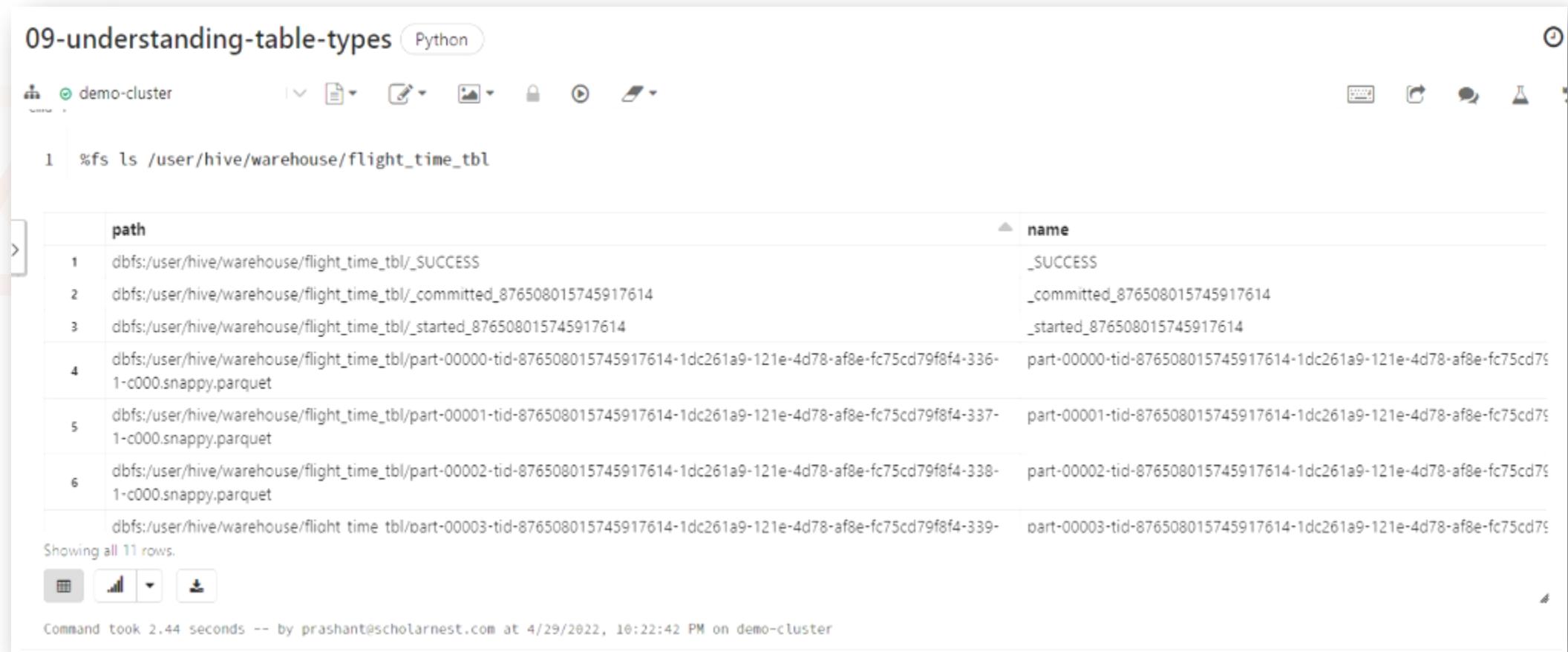
Now, moving on to the second topic.

Spark manages table metadata and table data together.

It means, Spark creates table data and metadata both when you create a managed table.

Similarly, Spark deletes the table data and metadata both when you drop the table.

We created the `flight_time_tbl`. So spark created data in the `/user/hive/warehouse/flight_time_tbl` directory. And Spark also created metadata in the metadata store. You can check the data directory using the `%fs ls` command as shown below. You will see the results, there are some parquet files here. These are your data files. Spark stored my data in the parquet files.



```
09-understanding-table-types Python

demo-cluster

1 %fs ls /user/hive/warehouse/flight_time_tbl
```

	path	name
1	dbfs:/user/hive/warehouse/flight_time_tbl/_SUCCESS	_SUCCESS
2	dbfs:/user/hive/warehouse/flight_time_tbl/_committed_876508015745917614	_committed_876508015745917614
3	dbfs:/user/hive/warehouse/flight_time_tbl/_started_876508015745917614	_started_876508015745917614
4	dbfs:/user/hive/warehouse/flight_time_tbl/part-00000-tid-876508015745917614-1dc261a9-121e-4d78-af8e-fc75cd79f8f4-336-1-c000.snappy.parquet	part-00000-tid-876508015745917614-1dc261a9-121e-4d78-af8e-fc75cd79f8f4-336-1-c000.snappy.parquet
5	dbfs:/user/hive/warehouse/flight_time_tbl/part-00001-tid-876508015745917614-1dc261a9-121e-4d78-af8e-fc75cd79f8f4-337-1-c000.snappy.parquet	part-00001-tid-876508015745917614-1dc261a9-121e-4d78-af8e-fc75cd79f8f4-337-1-c000.snappy.parquet
6	dbfs:/user/hive/warehouse/flight_time_tbl/part-00002-tid-876508015745917614-1dc261a9-121e-4d78-af8e-fc75cd79f8f4-338-1-c000.snappy.parquet	part-00002-tid-876508015745917614-1dc261a9-121e-4d78-af8e-fc75cd79f8f4-338-1-c000.snappy.parquet
	dbfs:/user/hive/warehouse/flight_time_tbl/part-00003-tid-876508015745917614-1dc261a9-121e-4d78-af8e-fc75cd79f8f4-339-	part-00003-tid-876508015745917614-1dc261a9-121e-4d78-af8e-fc75cd79f8f4-339-

Showing all 11 rows.

Command took 2.44 seconds -- by prashant@scholarnest.com at 4/29/2022, 10:22:42 PM on demo-cluster

We have two methods to query the metadata store:

1. SQL Commands
2. Spark Catalog API

The screenshot shows a Jupyter Notebook cell with the following content:

```
1 %sql
2 show tables
```

	database	tableName	isTemporary
1	default	flight_time_tbl	false

Showing all 1 rows.

Command took 0.39 seconds -- by prashant@scholarnest.com at 4/29/2022, 10:23:27 PM on demo-cluster

Cmd 9

```
1 spark.catalog.listTables()
```

▶ (1) Spark Jobs

Out[10]: [Table(name='flight\_time\_tbl', database='default', description=None, tableType='MANAGED', isTemporary=False)]

Command took 0.82 seconds -- by prashant@scholarnest.com at 4/29/2022, 10:24:42 PM on demo-cluster

Cmd 10

Now let me drop the table.

It was a managed table, so Spark will delete the table metadata and also delete the data files.

Cmd 10

```
1 %sql  
2 drop table flight_time_tbl
```

OK

Command took 4.03 seconds -- by prashant@scholarnest.com at 4/29/2022, 10:25:23 PM on demo-cluster

Now you can run your `%fs ls` command, and you will see an `FileNotFoundException`. So the table data is gone. Let me re-run the show tables command. And you can see the Query returned no results. So the metadata is also gone.

The screenshot shows a Jupyter Notebook cell with the title "09-understanding-table-types" and a Python tab selected. The cell contains the following code and output:

```
demo-cluster
Command took 0.87 seconds -- by prashant@scholarnest.com at 4/29/2022, 10:18:15 PM on demo-cluster
Cmd 7
> 1 %fs ls /user/hive/warehouse/flight_time_tbl
@FileNotFoundException: /user/hive/warehouse/flight_time_tbl
Command took 0.85 seconds -- by prashant@scholarnest.com at 4/29/2022, 10:26:08 PM on demo-cluster
Cmd 8
1 %sql
2 show tables

Query returned no results
Command took 0.15 seconds -- by prashant@scholarnest.com at 4/29/2022, 10:26:13 PM on demo-cluster
```

Now let's talk about the external tables.

The external tables are also referred to as unmanaged tables.

However, the external table is a more popular name.

So what are external tables?

Spark offers an external table so you can create a spark table over shared external data.

External tables are a mechanism for sharing data across projects or storage layers.

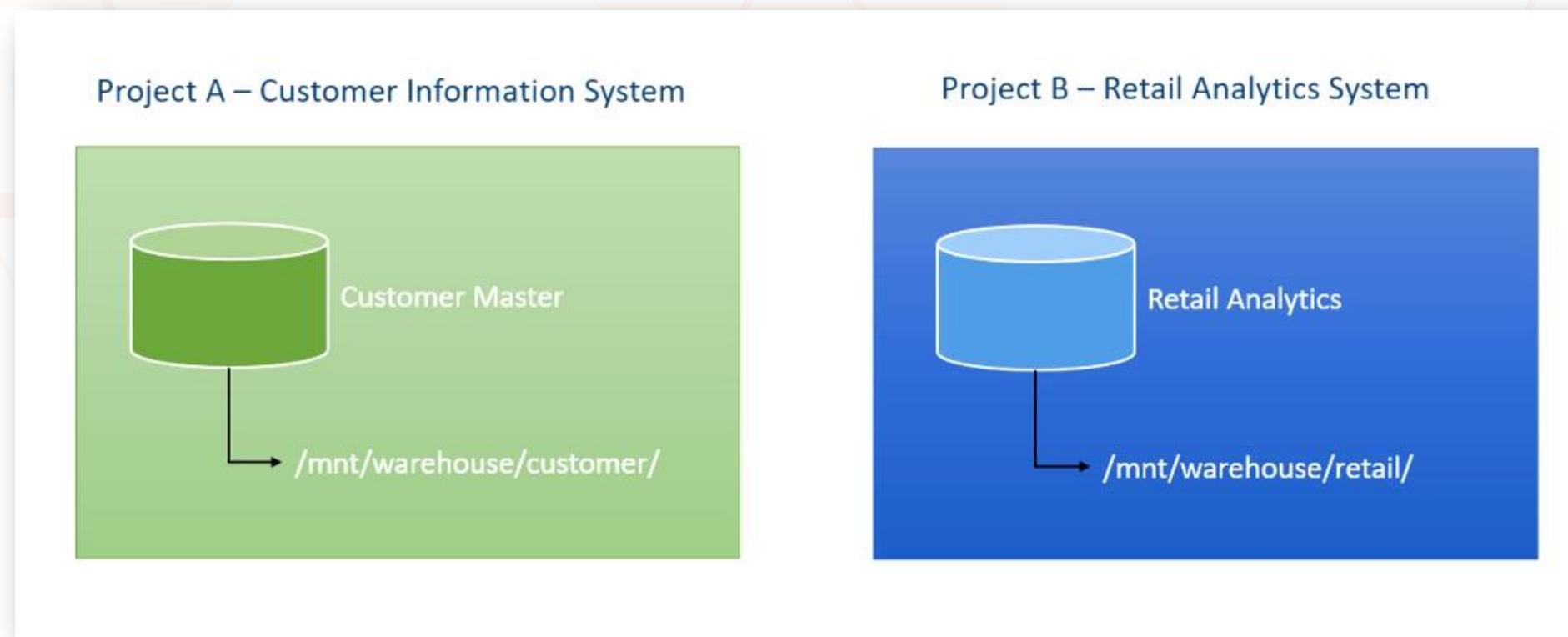
Assume you have two projects. The first project collects and manages customer information. So you have a customer master table in the project. The table data is stored in some directory in the project warehouse directory location. Now you are starting a new project for retail data analysis. This project has a different warehouse directory location where you store the data.

We have two separate projects and two different data stores.

Both the projects belong to the same organization, but we have two different data owners.

Project A owns the customer data, and project B owns the retail transactions.

But you have a requirement in Project B to access and use the customer master table.



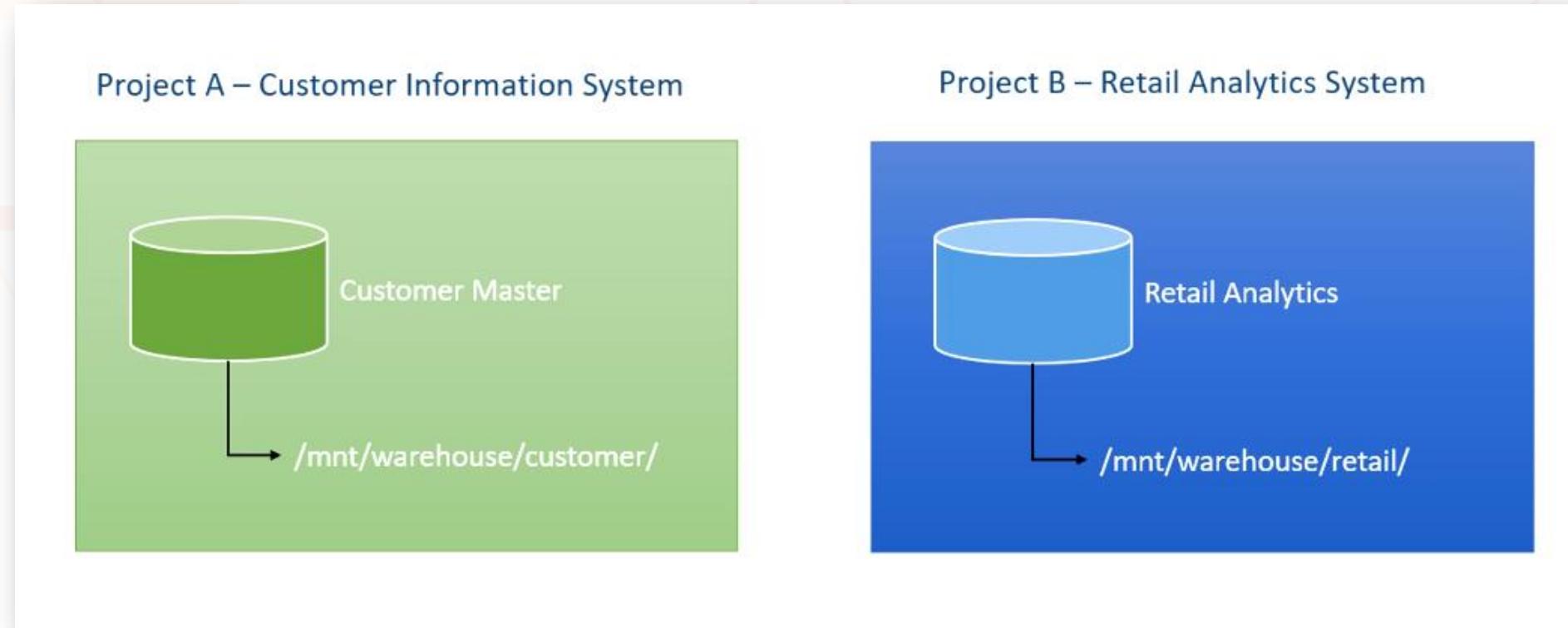
We have two approaches to solve our problem:

1. Copy the customer master data from project A to project B and use it.
2. Create an external table from project A to project B and use the same dataset.

The first approach is problematic.

Copying data at two places requires extra space, and you also need to keep both copies in sync.

So it is much better to use the second approach.

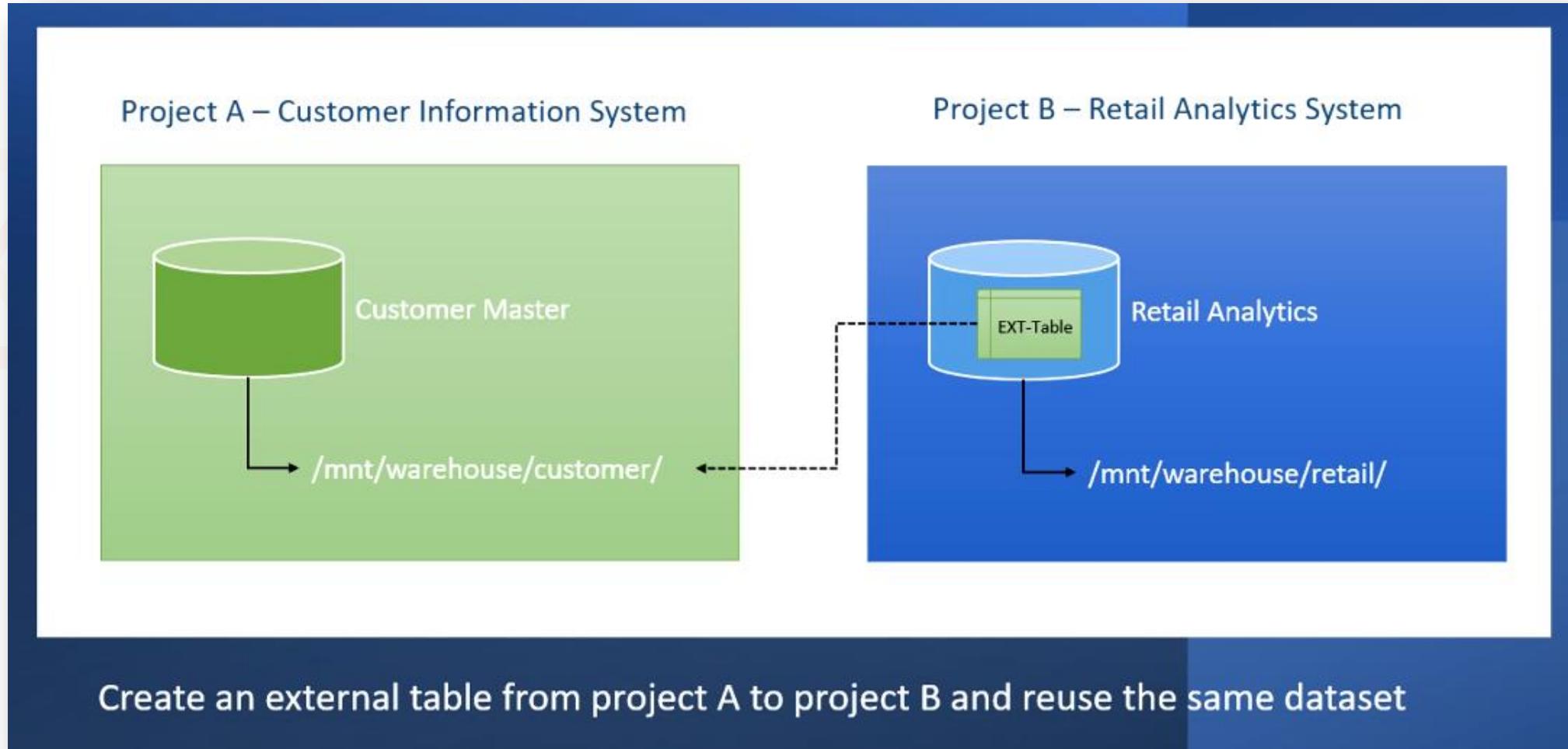


We can create an external table from project A to Project B and use the same data.

So the idea behind the external table is to share data across ownerships and storage layers.

Project A could be using a different storage layer than project B.

But they can share the same data via external tables as long as Spark supports both storage layers.



I dropped the flight\_time\_tbl table, so let me recreate it.

Assume the flight\_time\_tbl is a Project A managed table. And Project B wants to reuse my table.

So they can create an external table on my flight\_time\_tbl and read data from it.

How to create an external table? They can use create table DDL as shown below.

Cmd 11

```
1 %sql
2 create table if not exists ext_flight_time_tbl like flight_time_tbl
3 location "/user/hive/warehouse/flight_time_tbl"
```

OK

Command took 0.79 seconds -- by prashant@scholarnest.com at 4/30/2022, 6:26:10 PM on demo-cluster

What is the LIKE in the statement given below? You can use the LIKE clause to copy the schema from an existing table. So it will look into the column definitions of flight\_time\_tbl and create the same columns and data types for ext\_flight\_time\_tbl.

So the LIKE clause is a quick way to copy the structure of an existing table.

This statement will create a managed table. But we wanted to create an external table.

You can set the location, and spark will create an external table. We want to share the same data, so we should point this ext\_flight\_time\_tbl to the same data files where the flight\_time\_tbl stores its data.

Cmd 11

```
1 %sql
2 create table if not exists ext_flight_time_tbl like flight_time_tbl
3 location "/user/hive/warehouse/flight_time_tbl"
```

OK

Command took 0.79 seconds -- by prashant@scholarnest.com at 4/30/2022, 6:26:10 PM on demo-cluster

Now you can describe it and check the table type and location.

So it is an external table, and it points to the same location. So, it means we have two tables:

1. flight\_time\_tbl
2. ext\_flight\_time\_tbl

Cmd 12

```
1 %sql
2 describe extended ext_flight_time_tbl
```

	col_name	data_type	comment
21	Created Time	Sat Apr 30 12:56:11 UTC 2022	
22	Last Access	UNKNOWN	
23	Created By	Spark 3.2.1	
24	Type	EXTERNAL	
25	Provider	parquet	
26	Location	dbfs:/user/hive/warehouse/flight_time_tbl	
27	Serde Library	org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe	

Showing all 29 rows.

Command took 0.48 seconds -- by prashant@scholarnest.com at 4/30/2022, 6:26:41 PM on demo-cluster

You can even check the catalog API, we have two tables:

1. flight\_time\_tbl
2. ext\_flight\_time\_tbl

```
1 spark.catalog.listTables()

▶ (2) Spark Jobs
> Out[6]: [Table(name='ext_flight_time_tbl', database='default', description=None, tableType='EXTERNAL', isTemporary=False),
   Table(name='flight_time_tbl', database='default', description=None, tableType='MANAGED', isTemporary=False)]
Command took 1.12 seconds -- by prashant@scholarnest.com at 4/30/2022, 6:27:33 PM on demo-cluster
End 14
```

You can count the records in the ext\_flight\_time\_tbl. I created the table and haven't loaded the data in this external table. But you can see the record count. Because the data already exists at the given location, and the external table simply points to the same data. You can use the external table in the same way as a managed table. I mean, you can read it, update it, overwrite it or do whatever you want. But you cannot drop it. I mean, you can drop an external table, but Spark will delete the metadata only.



Cmd 14

```
1 %sql
2 select count(*) from ext_flight_time_tbl
```

▶ (2) Spark Jobs

count(1)
1 300000

Showing all 1 rows.

Command took 9.42 seconds -- by prashant@scholarnest.com at 4/30/2022, 6:28:22 PM on demo-cluster

Try the DROP statement on the external table. And now if you query the catalog, it will show you only one table.

The screenshot shows a Jupyter Notebook cell titled "09-understanding-table-types" in Python mode. The cell contains the following code:

```
1 %sql
2 drop table ext_flight_time_tbl
```

The output of the cell shows the command was successful:

> OK

Command took 0.73 seconds -- by prashant@scholarnest.com at 4/30/2022, 6:29:26 PM on demo-cluster

The cell then lists the tables in the catalog:

```
1 spark.catalog.listTables()
```

The output shows one table:

▶ (1) Spark Jobs

Out[7]: [Table(name='flight\_time\_tbl', database='default', description=None, tableType='MANAGED', isTemporary=False)]

Command took 0.66 seconds -- by prashant@scholarnest.com at 4/30/2022, 6:29:55 PM on demo-cluster

Cmd 17

So the external table is gone. But data is still there. And it makes sense also because the original table exists. If you query the original table as shown below, you can see the records here.

```
Cmd 17

1 %sql
2 select count(*) from flight_time_tbl

▶ (2) Spark Jobs

  count(1) ▲
  1  300000

Showing all 1 rows.

  □  ▢ ▾  ▴
```

Command took 2.50 seconds -- by prashant@scholarnest.com at 4/30/2022, 6:30:30 PM on demo-cluster

## Spark Tables

---

### 1. Managed Tables

- Created at a predefined warehouse directory location
- Drop Table will delete data and metadata both

### 2. External Tables

- Shared data across project/storage
- Set the location/path to create an external table
- Allows all operations such as append/overwrite
- Best practice is to avoid append/overwrite external tables
- Drop table will delete metadata only



Thank You  
ScholarNest Technologies Pvt Ltd.  
[www.scholarnest.com](http://www.scholarnest.com)



ScholarNest

# Spark Azure Databricks

Databricks Spark Certification and beyond

Instructor: Prashant Kumar Pandey



Absolute Beginner to Specialization in Apache Spark and Azure Databricks





# Spark Web UI - Introduction

## Spark Web UI

The go-to place to monitor the status and resource consumption of your Spark cluster

Go to your Databricks workspace.

Now, go to compute menu and start a cluster.

I already have one cluster running.

Click the cluster and go to the cluster details page.

The screenshot shows the Databricks Compute interface. On the left, there's a sidebar with icons for Compute, Databricks Community, and Help. The main area is titled 'Compute' and shows a table of clusters. The table has columns for Name, State, Nodes, Runtime, Driver, Worker, Creator, and Actions. One cluster is listed: 'demo-cluster' (State: Running, 1 node, 10.4 LTS runtime), created by 'prashant@sc...', and owned by 'Community'. There are buttons for 'Create Cluster' and 'Edit Cluster'.

Name	State	Nodes	Runtime	Driver	Worker	Creator	Actions
demo-cluster	Running	1 (0 spot)	10.4 LTS (includes Apache Spark 3.2.1, Scala 2....)	Community...	Community...	prashant@sc...	...

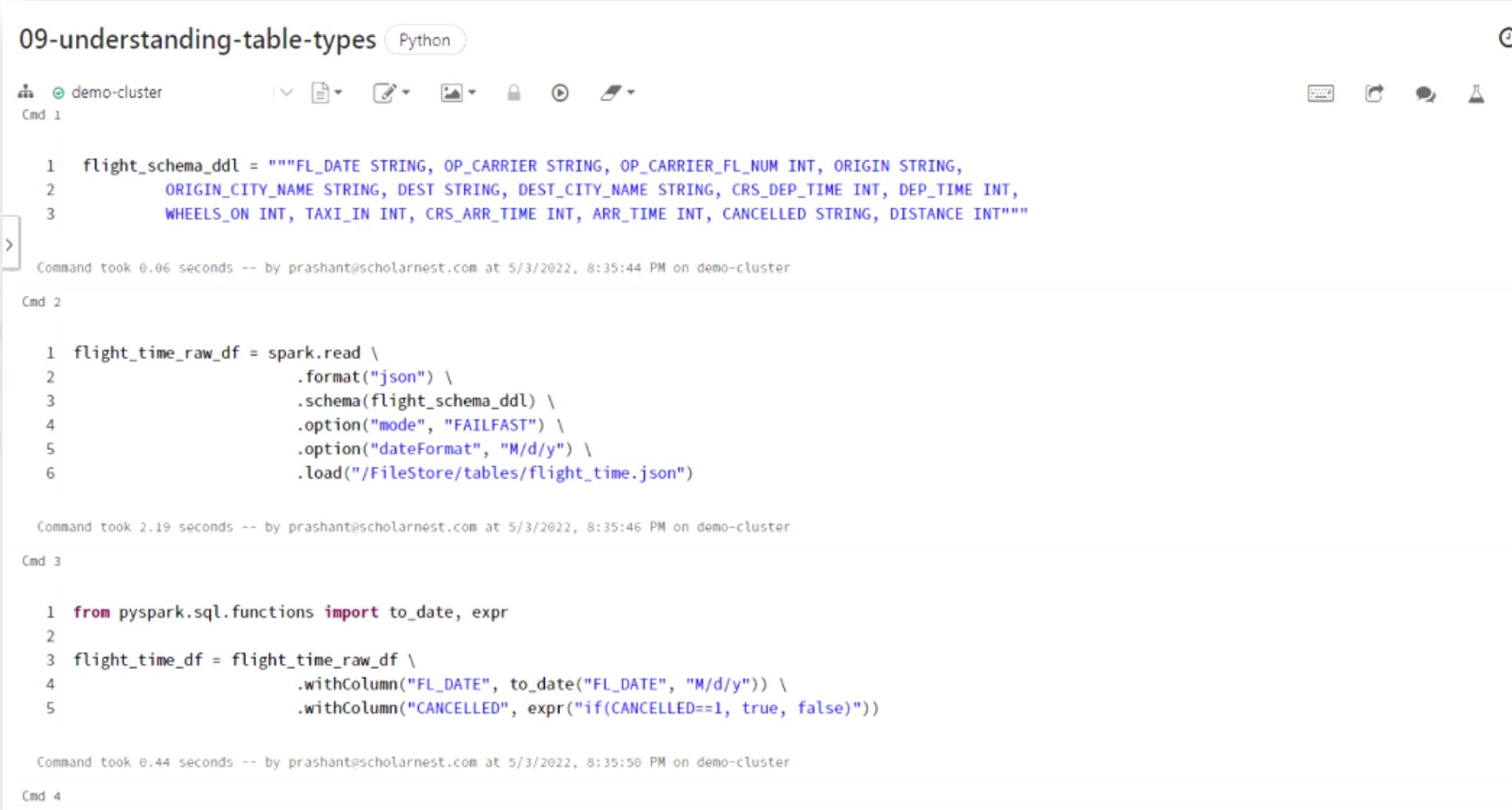
Then, navigate to your Spark Web UI. You will see some tabs such as Jobs, Stages, storage, etc. Spark UI is your go-to place for monitoring your Spark applications and understanding many other details about your application.

You do not see anything here because I haven't run anything on this cluster.

Things will start appearing here as you start running your application or the Spark code.

The screenshot shows the Spark Web UI interface for a cluster named "demo-cluster". The top navigation bar includes tabs for Configuration, Notebooks, Libraries, Event log, Spark UI (which is highlighted with a yellow box), Driver logs, Metrics, Apps, and Spark cluster UI - Master. Below the tabs, there are several status indicators: Hostname (ec2-35-163-128-95.us-west-2.compute.amazonaws.com), Spark Version (10.4.x-scala2.12), and a user summary (User: root, Total Uptime: 7.0 min, Scheduling Mode: FAIR). A red box highlights the "Jobs" tab under the sub-menu "Spark Jobs". At the bottom left, there is a link to "Event Timeline".

Open a notebook from the previous lecture. (**Reference - 09-understanding-table-types**)  
Attach a cluster to this notebook. And run all the cells up to where we are writing a Dataframe.



09-understanding-table-types Python

demo-cluster

Cmd 1

```
1 flight_schema_ddl = """FL_DATE STRING, OP_CARRIER STRING, OP_CARRIER_FL_NUM INT, ORIGIN STRING,  
2 ORIGIN_CITY_NAME STRING, DEST STRING, DEST_CITY_NAME STRING, CRS_DEP_TIME INT, DEP_TIME INT,  
3 WHEELS_ON INT, TAXI_IN INT, CRS_ARR_TIME INT, ARR_TIME INT, CANCELLED STRING, DISTANCE INT"""
```

Command took 0.06 seconds -- by prashant@scholarnest.com at 5/3/2022, 8:35:44 PM on demo-cluster

Cmd 2

```
1 flight_time_raw_df = spark.read \  
2 .format("json") \  
3 .schema(flight_schema_ddl) \  
4 .option("mode", "FAILFAST") \  
5 .option("dateFormat", "M/d/y") \  
6 .load("/FileStore/tables/flight_time.json")
```

Command took 2.19 seconds -- by prashant@scholarnest.com at 5/3/2022, 8:35:46 PM on demo-cluster

Cmd 3

```
1 from pyspark.sql.functions import to_date, expr  
2  
3 flight_time_df = flight_time_raw_df \  
4 .withColumn("FL_DATE", to_date("FL_DATE", "M/d/y")) \  
5 .withColumn("CANCELLED", expr("if(CANCELLED==1, true, false)"))
```

Command took 0.44 seconds -- by prashant@scholarnest.com at 5/3/2022, 8:35:50 PM on demo-cluster

Cmd 4

Click the demo-cluster drop down, and you will see a direct link for the Spark UI.

Let me go there.

The screenshot shows a Jupyter Notebook interface with the following details:

- Title:** 09-understanding-table-types
- Languages:** Python
- Attached cluster:** demo-cluster (selected)
- Cluster Details:** 15.25 GB | 2 Cores | DMR 10.4 LTS | Spark 3.2.1 | Scala 2.12
- Action Buttons:** Detach, Restart Cluster, Detach & Re-attach, **Spark UI** (highlighted with a red box), Driver logs, Terminal
- Content Area:** (1) Spark Jobs
- Log Message:** Command took 36.48 seconds -- by prashant@scholarnest.com at 5/3/2022, 8:36:08 PM on demo-cluster

You can see some information here. We are currently in the Jobs tab of the Spark UI. The Jobs tab displays a summary page of all jobs in the Spark application. The summary page shows high-level information, such as the status, duration, and progress bar.

The screenshot shows the Databricks Shell - Spark Jobs interface. At the top, there are tabs for Jobs, Stages, Storage, Environment, Executors, SQL, JDBC/ODBC Server, and Structured Streaming. The Jobs tab is selected. Below the tabs, it says "User: root", "Total Uptime: 11 min", "Scheduling Mode: FAIR", and "Completed Jobs: 1". There are links for "Event Timeline" and "Completed Jobs (1)". The "Completed Jobs (1)" section is expanded, showing a table with one row. The table has columns: Job Id (Job Group), Description, Submitted, Duration, Stages: Succeeded/Total, and Tasks (for all stages): Succeeded/Total. The row details a job with ID 0, submitted on 2022/05/03 at 15:06:10, taking 30 s, with 1/1 stages succeeded and 8/8 tasks succeeded. Navigation controls at the bottom allow for page selection (Page: 1), jumping to a specific page (1 Pages. Jump to 1), and showing 100 items per page (Show 100 items in a page, Go).

Job Id (Job Group) *	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
0 (1901202411884247856_5170693221701987476_a5f884691dec46bd892089b3a869635e)	flight_time_df.write \.format("par... saveAsTable at NativeMethodAccessorImpl.java:0	2022/05/03 15:06:10	30 s	1/1	8/8

You can click on a job description on the summary page to see the details.

Jobs Stages Storage Environment Executors SQL JDBC/ODBC Server Structured Streaming

Spark Jobs (?)

User: root  
Total Uptime: 11 min  
Scheduling Mode: FAIR  
Completed Jobs: 1

► Event Timeline  
▼Completed Jobs (1)

Page: 1 1 Pages. Jump to 1 . Show 100 items in a page. Go

Job Id (Job Group) ▾	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
0 (1901202411884247856_5170693221701987476_a5f884691dec46bd892089b3a869635e)	flight_time_df.write \ .format("par... saveAsTable at NativeMethodAccessorImpl.java:0	2022/05/03 15:06:10	30 s	1/1	8/8

Page: 1 1 Pages. Jump to 1 . Show 100 items in a page. Go



Here we are on the Job details for job id zero.

The Job details page shows you the stage-wise breakup of your Job. I have only one stage in this Job. The details page also shows some other information about the Job stages, such as duration and task progress bar.

The screenshot shows the 'Jobs' tab selected in the navigation bar. Below it, the title 'Details for Job 0' is displayed. Key metrics shown include:

- Status:** SUCCEEDED
- Submitted:** 2022/05/03 15:06:10
- Duration:** 30 s
- Associated SQL Query:** 4
- Job Group:** 1901202411884247856\_5170693221701987476\_a5f884691dec46bd892089b3a869635e
- Completed Stages:** 1

Below these details, there are links for 'Event Timeline' and 'DAG Visualization'. A section titled 'Completed Stages (1)' is expanded, showing a table with the following data:

Stage Id	Pool Name	Description	Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write
0	1901202411884247856	flight_time_df.write \.format("par... saveAsTable at NativeMethodAccessorImpl.java:0	+details 2022/05/03 15:06:11	30 s	8/8	3.6 MiB			

At the bottom of the page, there are two pagination sections, each with 'Page: 1' and '1 Pages. Jump to 1 . Show 100 items in a page. Go' buttons.

Click the DAG visualization link to see the DAG for the selected Job.

You might not understand these things yet because I didn't cover Spark Jobs, Stages, and Tasks. However, the information shown over the Spark UI is crucial for understanding your application behavior and resources and identifying tuning and optimization opportunities.



I have only one Job as of now. It makes sense to run a few more jobs to understand how this page looked when we ran multiple jobs on the same cluster. So come back to your notebook and run a few more commands. Run the create external table command, come down to the select count(\*) from ext\_flight\_time\_tbl, and run that as well.

Now, look at the output. It says 2 Spark Jobs and also shows the count. So the select statement triggered two Spark jobs.

Cmd 14

```
1 %sql
2 select count(*) from ext_flight_time_tbl
```

▶ (2) Spark Jobs ←

	count(1)
1	300000

Showing all 1 rows.

Command took 9.85 seconds -- by prashant@scholarnest.com at 5/3/2022, 8:42:44 PM on demo-cluster

Go back to the Spark UI once again and refresh the page, and you will see three jobs. We can see Job IDs in the first column. Job-0 was triggered by the `saveAsTable()` method. The `saveAsTable()` is a `DataFramewriter` method. Job-1 and Job-2 were triggered by the `select count(*) SQL`.

The screenshot shows the Databricks Spark Jobs UI. At the top, it displays the URL `community.cloud.databricks.com/sparkui/0503-145506-vzyghp07/driver-6906967366313043515/jobs/?o=946780288864704`. Below the header, there are tabs for Jobs, Stages, Storage, Environment, Executors, SQL, JDBC/ODBC Server, and Structured Streaming. The Jobs tab is selected.

The main area is titled "Spark Jobs (?)". It shows the following information:

- User: root
- Total Uptime: 18 min
- Scheduling Mode: FAIR
- Completed Jobs: 3

Below this, there are two sections: "Event Timeline" and "Completed Jobs (3)". The "Completed Jobs (3)" section is expanded, showing three rows of job details:

Job Id (Job Group)	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
2 (3094848351750946294_5576904425099990681_042e26e8d430464e88e38248396474c1)	select count(*) from ext_flight_time_tbl collectResult at OutputAggregator.scala:268	2022/05/03 15:12:53	0.3 s	1/1 (1 skipped)	1/1 (8 skipped)
1 (3094848351750946294_5576904425099990681_042e26e8d430464e88e38248396474c1)	select count(*) from ext_flight_time_tbl collectResult at OutputAggregator.scala:268	2022/05/03 15:12:45	7 s	1/1	8/8
0 (1901202411884247856_5170693221701987476_a5f884691dec46bd892089b3a869635e)	flight_time_df.write \.format("par... saveAsTable at NativeMethodAccesso...impl.java:0	2022/05/03 15:06:10	30 s	1/1	8/8

Now, go to the next tab that shows the stage summary page.

So the stage summary page shows all the stages across the Jobs.

We ran three jobs, and we saw them on the Jobs summary page.

Altogether, those four jobs ran in four stages, and we can see four stages listed here on this stages summary page.

Jobs	Stages	Storage	Environment	Executors	SQL	JDBC/ODBC Server	Structured Streaming																
Completed Stages: 3																							
Skipped Stages: 1																							
▼ Fair Scheduler Pools (1)																							
<table><thead><tr><th>Pool Name</th><th>Minimum Share</th><th>Pool Weight</th><th>Active Stages</th><th>Running Tasks</th><th>SchedulingMode</th><th> </th><th> </th></tr></thead><tbody><tr><td>default</td><td>0</td><td>1</td><td>0</td><td>0</td><td>FIFO</td><td> </td><td> </td></tr></tbody></table>								Pool Name	Minimum Share	Pool Weight	Active Stages	Running Tasks	SchedulingMode			default	0	1	0	0	FIFO		
Pool Name	Minimum Share	Pool Weight	Active Stages	Running Tasks	SchedulingMode																		
default	0	1	0	0	FIFO																		
▼ Completed Stages (3)																							
Page: 1				1 Pages. Jump to <input type="text" value="1"/> . Show <input type="text" value="100"/> items in a page. <input type="button" value="Go"/>																			
Stage Id	Pool Name	Description	Submitted	Duration	Tasks: Succeeded/Total	Input	Output																
3	3094848351750946294	select count(*) from ext_flight_time_tbl collectResult at OutputAggregator.scala:268	2022/05/03 15:12:53	0.2 s	1/1	472.0 B																	
1	3094848351750946294	select count(*) from ext_flight_time_tbl collectResult at OutputAggregator.scala:268	2022/05/03 15:12:46	7 s	8/8		472.0 B																
0	1901202411884247856	flight_time_df.write \.format("par... saveAsTable at NativeMethodAccessorImpl.java:0	2022/05/03 15:06:11	30 s	8/8	3.6 MiB																	
Page: 1				1 Pages. Jump to <input type="text" value="1"/> . Show <input type="text" value="100"/> items in a page. <input type="button" value="Go"/>																			
▼ Skipped Stages (1)																							
Page: 1				1 Pages. Jump to <input type="text" value="1"/> . Show <input type="text" value="100"/> items in a page. <input type="button" value="Go"/>																			
Stage Id	Pool Name	Description	Submitted	Duration	Tasks: Succeeded/Total	Input	Output																
2	default	collectResult at OutputAggregator.scala:268	+details	Unknown	Unknown	0/8																	

Go to the Jobs tab.

I have three jobs, and I know Job-0 is saving my Dataframe to a Spark table.

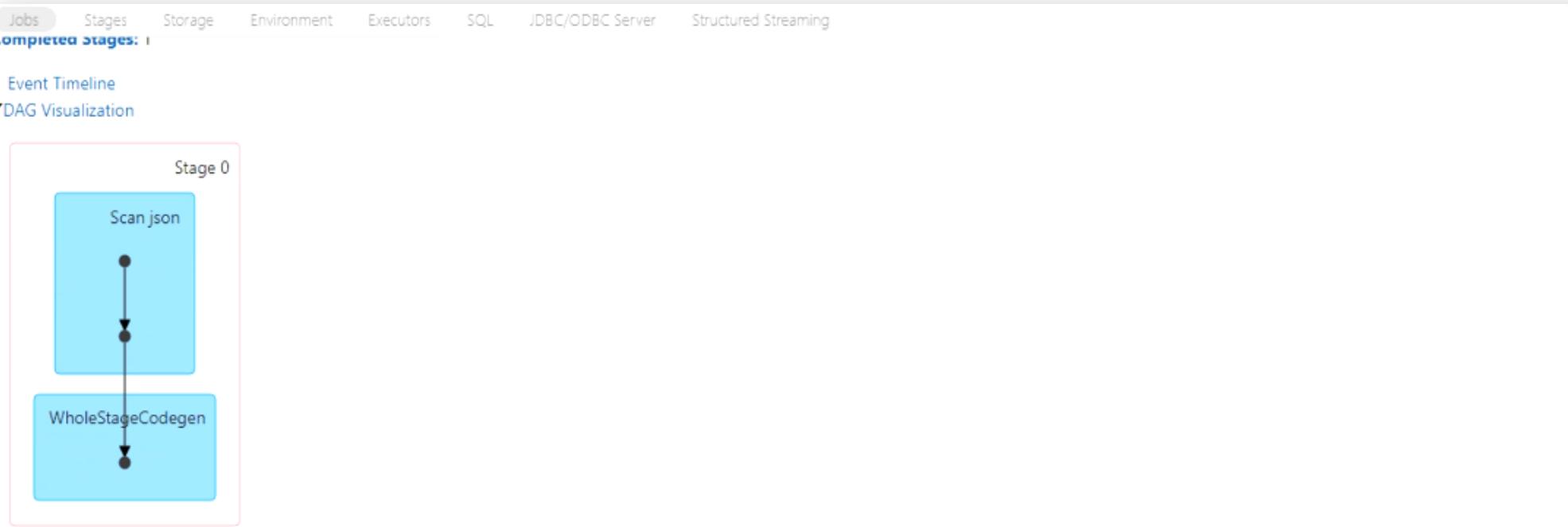
I also know that Job-1 and Job-2 give me counts from the external table.

The Job-0 took 30 seconds, so I want to investigate it further.

So I will click the Job description and reach the Job details page.

Spark Jobs (?)						
User: root						
Total Uptime: 43 min						
Scheduling Mode: FAIR						
Completed Jobs: 3						
<a href="#">Event Timeline</a>						
<a href="#">Completed Jobs (3)</a>						
Page: 1		1 Pages. Jump to <input type="text" value="1"/> . Show <input type="text" value="100"/> items in a page. <a href="#">Go</a>				
Job Id (Job Group) *		Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
2 (3094848351750946294_5576904425099990681_042e26e8d430464e88e38248396474c1)		select count(*) from ext_flight_time_tbl collectResult at OutputAggregator.scala:268	2022/05/03 15:12:53	0.3 s	1/1 (1 skipped)	<span style="background-color: #00aaff; color: white; padding: 2px;">1/1 (8 skipped)</span>
1 (3094848351750946294_5576904425099990681_042e26e8d430464e88e38248396474c1)		select count(*) from ext_flight_time_tbl collectResult at OutputAggregator.scala:268	2022/05/03 15:12:45	7 s	1/1	<span style="background-color: #00aaff; color: white; padding: 2px;">8/8</span>
0 (1901202411884247856_5170693221701987476_a5f884691dec46bd892089b3a869635e)		flight_time_df.write \.format("par... saveAsTable at NativeMethodAccessorImpl.java:0	2022/05/03 15:06:10	30 s	1/1	<span style="background-color: #00aaff; color: white; padding: 2px;">8/8</span>

I see only one stage here, and I know this stage belongs to Job-0.  
I want to look into the stage details, so I will click the Stage description and go to the stage details page.



The DAG visualization shows Stage 0 with two tasks: "Scan json" and "WholeStageCodegen". The "Scan json" task is at the top, followed by a vertical line with two dots, leading to the "WholeStageCodegen" task at the bottom. Both tasks are represented by blue rectangular boxes with black outlines.

**Completed Stages (1)**

Stage Id	Pool Name	Description	Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write
0	1901202411884247856	flight_time_df.write \.format("par... saveAsTable at NativeMethodAccessorImpl.java:0	+details	2022/05/03 15:06:11	30 s	8/8	3.6 MiB		

Page: 1      1 Pages. Jump to 1 . Show 100 items in a page. Go

Page: 1      1 Pages. Jump to 1 . Show 100 items in a page. Go

Here I am in the stage details page. I can see the stage details of my Job-0. You can come here from the stage summary page also.

Jobs Stages Storage Environment Executors SQL JDBC/ODBC Server Structured Streaming

### Details for Stage 0 (Attempt 0)

**Resource Profile Id:** 0  
**Total Time Across All Tasks:** 3.8 min  
**Locality Level Summary:** Process local: 8  
**Output Size / Records:** 3.6 MiB / 300000  
**Associated Job Ids:** 0

▶ DAG Visualization  
▶ Show Additional Metrics  
▶ Event Timeline

#### Summary Metrics for 8 Completed Tasks

Metric	Min	25th percentile	Median	75th percentile	Max
Duration	28 s	28 s	28 s	29 s	29 s
GC Time	2 s	2 s	2 s	2 s	2 s
Output Size / Records	316.1 KiB / 25090	473.7 KiB / 39257	480.7 KiB / 39276	486.7 KiB / 39351	489.3 KiB / 39379

Showing 1 to 3 of 3 entries

▶ Aggregated Metrics by Executor

#### Tasks (8)

Show 20 entries Search:

Index	Task ID	Attempt	Status	Locality level	Executor ID	Host	Logs	Launch Time	Duration	GC Time	Output Size / Records	Errors
0	0	0	SUCCESS	PROCESS_LOCAL	driver	ip-10-172-228-131.us-west-2.compute.internal		2022-05-04	29 s	2 s	489.3 KiB / 39379	0

We also have a storage tab.

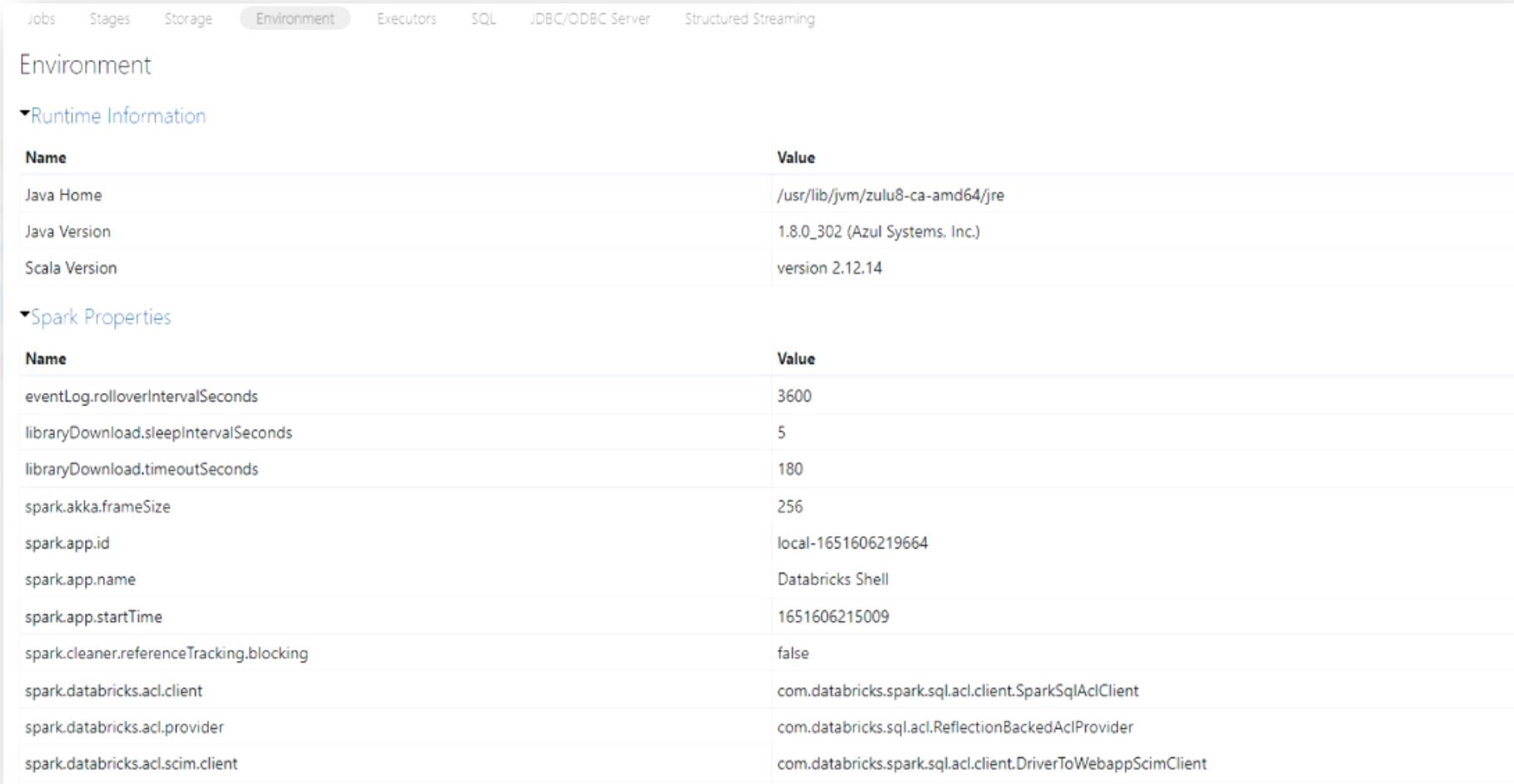
We do not see enough information here, but the storage page is your gateway to see the cache information.

This page shows information when you cache something in Spark memory

The screenshot shows the Apache Spark Web UI with the 'Storage' tab selected. The main title is 'Storage' and the sub-section is 'Parquet IO Cache'. Below this, there is a detailed table with the following data:

Data Read from External Filesystem (All Formats)	Data Read from IO Cache (Cache Hits, Compressed)	Data Written to IO Cache (Compressed)	Cache Misses (Compressed)	True Cache Misses	Partial Cache Misses	Rescheduling Cache Misses	Cache Hit Ratio	Number of Local Scan Tasks	Number of Rescheduled Scan Tasks	Cache Metadata Manager Peak Disk Usage
83.6 MiB	0.0 B	0.0 B	0.0 B	0.0 B	0.0 B	0.0 B	0 %	0	0	0.0 B

The next page is for the environment.  
The Environment tab displays the values for the different environment and configuration variables, including JVM, Spark, and system properties.  
We refer to this page to understand the runtime environment and see if some tuning is required.



The screenshot shows the Databricks UI with the 'Environment' tab selected. The page is titled 'Environment' and contains two sections: 'Runtime Information' and 'Spark Properties'. Both sections are presented as tables with 'Name' and 'Value' columns.

Name	Value
Java Home	/usr/lib/jvm/zulu8-ca-amd64/jre
Java Version	1.8.0_302 (Azul Systems, Inc.)
Scala Version	version 2.12.14

Name	Value
eventLog.rolloverIntervalSeconds	3600
libraryDownload.sleepIntervalSeconds	5
libraryDownload.timeoutSeconds	180
spark.akka.frameSize	256
spark.app.id	local-1651606219664
spark.app.name	Databricks Shell
spark.app.startTime	1651606215009
spark.cleaner.referenceTracking.blocking	false
spark.databricks.acl.client	com.databricks.spark.sql.acl.client.SparkSqlAclClient
spark.databricks.acl.provider	com.databricks.sql.acl.ReflectionBackedAclProvider
spark.databricks.acl.scim.client	com.databricks.spark.sql.acl.client.DriverToWebappScimClient

The next tab is the executor tab.

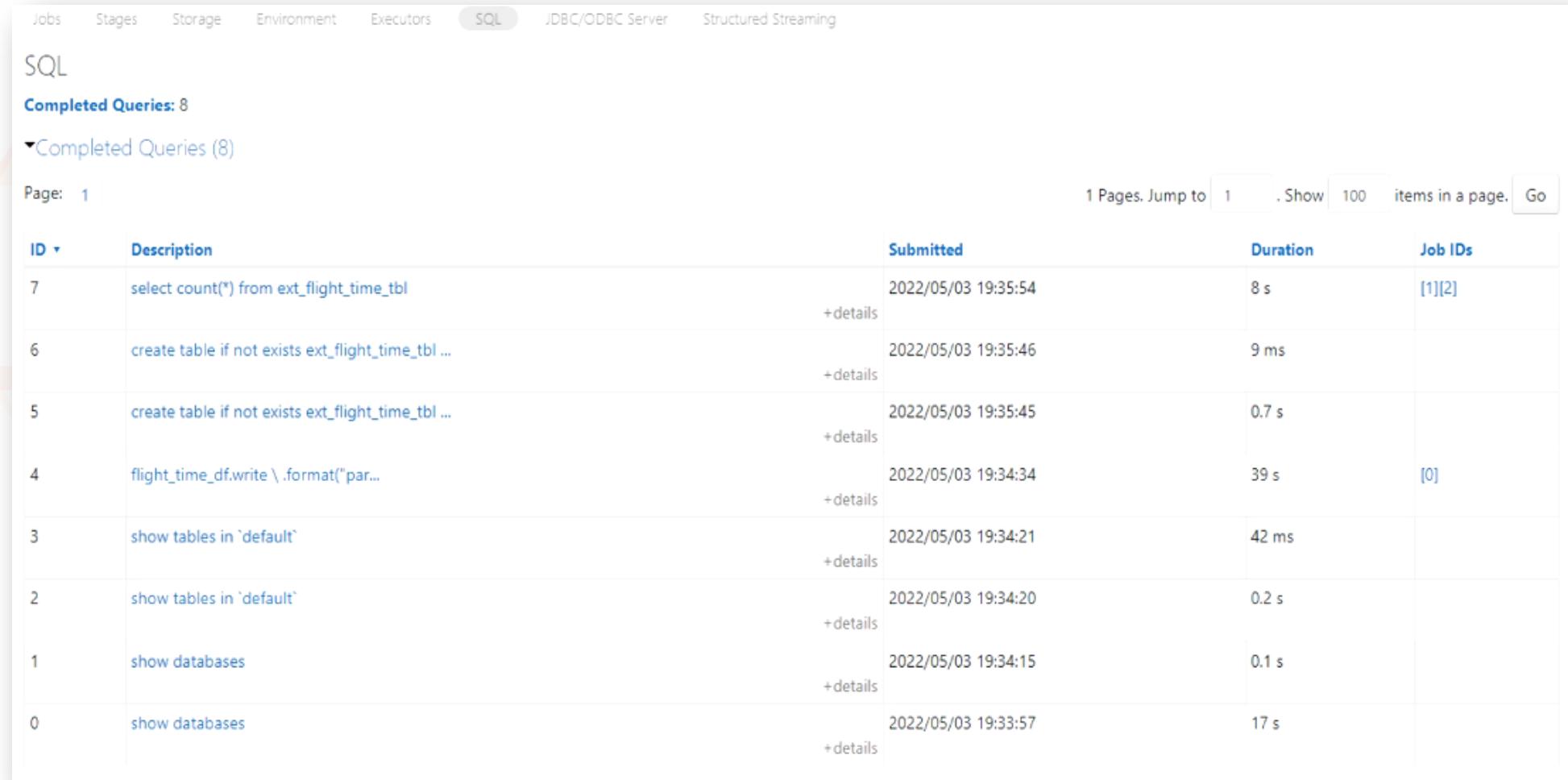
This guy displays summary information about the executors, including memory, cores, disk usage, and shuffle information each executor uses.

Jobs	Stages	Storage	Environment	Executors	SQL	JDBC/ODBC Server	Structured Streaming										
Executors																	
▶ Show Additional Metrics																	
Summary																	
▲	RDD Blocks	Storage Memory	Disk Used	Cores	Active Tasks	Failed Tasks	Complete Tasks	Total Tasks	Task Time (GC Time)	Input	Shuffle Read	Shuffle Write	Excluded				
Active(1)	0	260.5 KiB / 3.9 GiB	0.0 B	8	0	0	17	17	5.1 min (17 s)	0.0 B	472 B	472 B	0				
Dead(0)	0	0.0 B / 0.0 B	0.0 B	0	0	0	0	0	0.0 ms (0.0 ms)	0.0 B	0.0 B	0.0 B	0				
Total(1)	0	260.5 KiB / 3.9 GiB	0.0 B	8	0	0	17	17	5.1 min (17 s)	0.0 B	472 B	472 B	0				
Executors								Search:									
Show	20	entries															
Executor ID	▲ Address	Status	RDD Blocks	Storage Memory	Disk Used	Cores	Active Tasks	Failed Tasks	Complete Tasks	Total Tasks	Task Time (GC Time)	Input	Shuffle Read	Shuffle Write	Thread Dump	Heap Histogram	Exec Loss Reason
driver	10.172.241.6:38291	Active	0	260.5 KiB / 3.9 GiB	0.0 B	8	0	0	17	17	5.1 min (17 s)	0.0 B	472 B	472 B	Thread Dump	Heap Histogram	
Showing 1 to 1 of 1 entries								Previous							Next		

The next page is for SQL.

This page shows the list of SQL statements executed on your cluster.

You can see the time taken by each SQL. It is a summary page, and you go to the SQL details page, clicking any SQL.



The screenshot shows the Apache Spark UI's SQL history page. At the top, there are tabs for Jobs, Stages, Storage, Environment, Executors, SQL (which is selected), JDBC/ODBC Server, and Structured Streaming. Below the tabs, it says "Completed Queries: 8". A dropdown menu "Completed Queries (8)" is open. On the left, there are navigation links for "Page: 1" and "1 Pages. Jump to 1 . Show 100 items in a page. Go". The main area is a table with the following data:

ID	Description	Submitted	Duration	Job IDs
7	<a href="#">select count(*) from ext_flight_time_tbl</a>	2022/05/03 19:35:54 +details	8 s	[1][2]
6	<a href="#">create table if not exists ext_flight_time_tbl ...</a>	2022/05/03 19:35:46 +details	9 ms	
5	<a href="#">create table if not exists ext_flight_time_tbl ...</a>	2022/05/03 19:35:45 +details	0.7 s	
4	<a href="#">flight_time_df.write \.format("par...</a>	2022/05/03 19:34:34 +details	39 s	[0]
3	<a href="#">show tables in `default`</a>	2022/05/03 19:34:21 +details	42 ms	
2	<a href="#">show tables in `default`</a>	2022/05/03 19:34:20 +details	0.2 s	
1	<a href="#">show databases</a>	2022/05/03 19:34:15 +details	0.1 s	
0	<a href="#">show databases</a>	2022/05/03 19:33:57 +details	17 s	

You also have JDBC/ODBC server tab.

This tab is visible only when your Spark cluster also runs a Thrift server and allows JDBC/ODBC connections.

We do not see anything here because we do not have any connections.

The last tab is for Spark Structured Streaming.

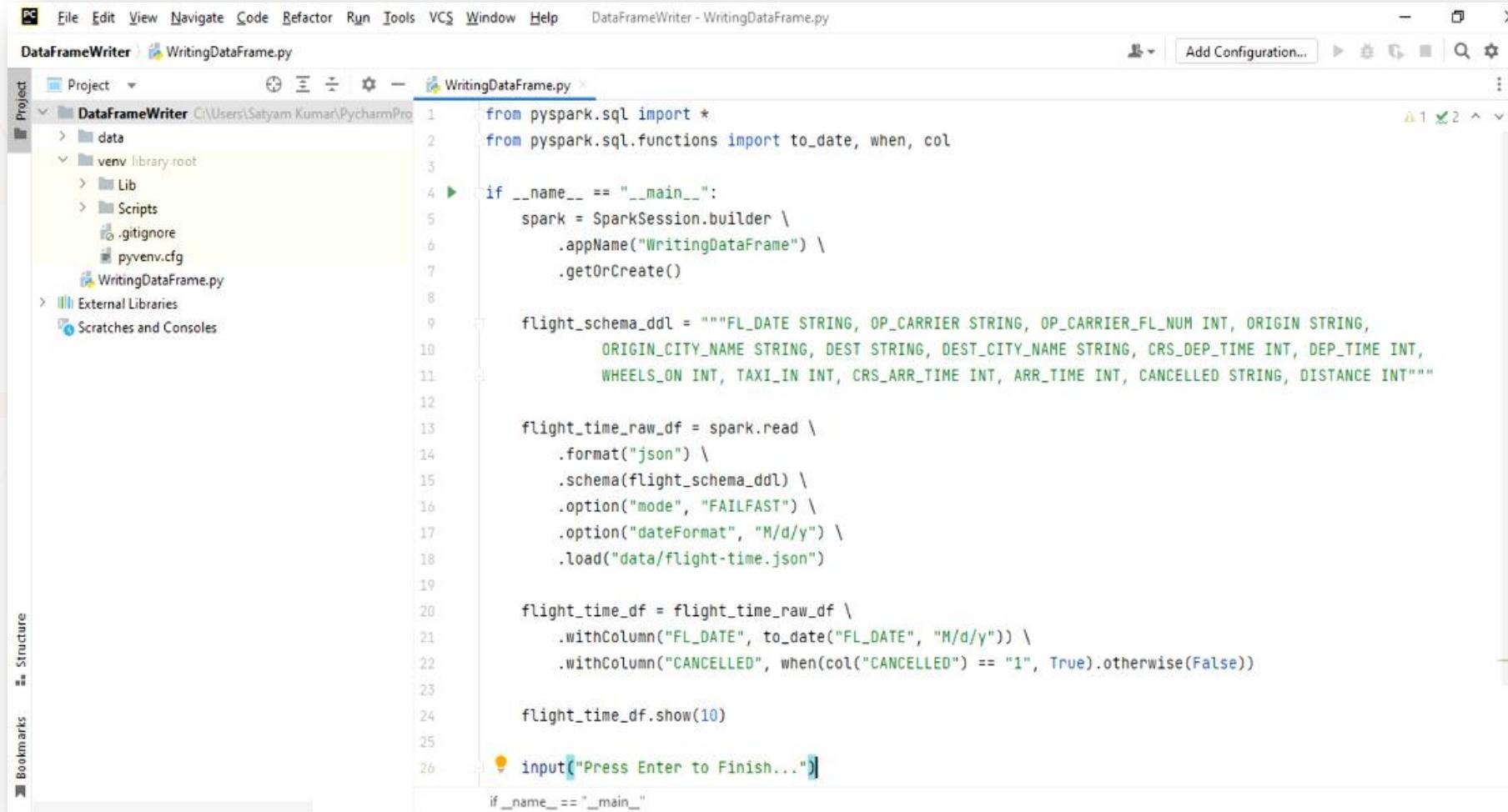
This tab shows information about streaming queries and micro-batches.

We will cover this tab in the Spark streaming course.

The screenshot shows the Apache Spark Web UI interface. At the top, there is a navigation bar with tabs: Jobs, Stages, Storage, Environment, Executors, SQL, JDBC/ODBC Server (which is highlighted), and Structured Streaming. Below the navigation bar, the main content area has a title "JDBC/ODBC Server". Underneath the title, it says "Started at: 2022/05/03 19:30:15" and "Time since start: 10 minutes 28 seconds". A message indicates "0 session(s) are online, running 0 SQL statement(s)". There are two expandable sections: "Session Statistics (0)" which states "No statistics have been generated yet.", and "SQL Statistics (0)" which also states "No statistics have been generated yet.".

You create and test the Spark application on your local machine. And you can still access the Spark UI. Here is my PyCharm IDE. I created a new project and wrote some code here. ([Reference - DataFrameWriter](#))

It is the same code that we saw in the notebook.



The screenshot shows the PyCharm IDE interface. The top menu bar includes File, Edit, View, Navigate, Code, Refactor, Run, Tools, VCS, Window, Help, and the current file name, DataFrameWriter - WritingDataFrame.py. The main window has a Project tool window on the left showing a directory structure for a project named 'DataFrameWriter' located at C:\Users\Satyam Kumar\PycharmPro. Inside the project are folders 'data', 'venv', 'Lib', 'Scripts', '.gitignore', 'pyenv.cfg', and the file 'WritingDataFrame.py'. The right side of the interface is the code editor with the following Python code:

```
PC File Edit View Navigate Code Refactor Run Tools VCS Window Help DataFrameWriter - WritingDataFrame.py
DataFrameWriter WritingDataFrame.py
Project Project - WritingDataFrame.py
DataFrameWriter C:\Users\Satyam Kumar\PycharmPro
> data
> venv library root
> Lib
> Scripts
  .gitignore
  pyenv.cfg
> WritingDataFrame.py
> External Libraries
Scratches and Consoles

from pyspark.sql import *
from pyspark.sql.functions import to_date, when, col
if __name__ == "__main__":
    spark = SparkSession.builder \
        .appName("WritingDataFrame") \
        .getOrCreate()

    flight_schema_ddl = """FL_DATE STRING, OP_CARRIER STRING, OP_CARRIER_FL_NUM INT, ORIGIN STRING,
                         ORIGIN_CITY_NAME STRING, DEST STRING, DEST_CITY_NAME STRING, CRS_DEP_TIME INT, DEP_TIME INT,
                         WHEELS_ON INT, TAXI_IN INT, CRS_ARR_TIME INT, ARR_TIME INT, CANCELLED STRING, DISTANCE INT"""

    flight_time_raw_df = spark.read \
        .format("json") \
        .schema(flight_schema_ddl) \
        .option("mode", "FAILFAST") \
        .option("dateFormat", "M/d/y") \
        .load("data/flight-time.json")

    flight_time_df = flight_time_raw_df \
        .withColumn("FL_DATE", to_date("FL_DATE", "M/d/y")) \
        .withColumn("CANCELLED", when(col("CANCELLED") == "1", True).otherwise(False))

    flight_time_df.show(10)

    input("Press Enter to Finish...")

if __name__ == "__main__"
```

You can run this application, and it will work. I added a Python input method, so my application is not terminating, and it is waiting for me. You can see in the output, there is a message “press enter to finish” don’t press enter yet. I am holding my application from termination to see the Spark UI.

```
File Edit View Navigate Code Refactor Run Tools VCS Window Help DataFrameWriter - WritingDataFrame.py
DataFrameWriter WritingDataFrame.py
Project  WritingDataFrame.py
DataFrameWriter C:\Users\Satyam Kumar\PycharmPro...
> data
> venv library root
> Lib
> Scripts
> .gitignore
> pyvenv.cfg
WritingDataFrame.py
23
24     flight_time_df.show(10)
25
26     input("Press Enter to Finish...")
if __name__ == "__main__":
22/05/04 11:49:47 WARN ProcfsMetricsGetter: Exception when trying to compute pagesize, as a result reporting of ProcessTree metrics is stopped
+-----+
| FL_DATE|OP_CARRIER|OP_CARRIER_FL_NUM|ORIGIN|ORIGIN_CITY_NAME|DEST|DEST_CITY_NAME|CRS_DEP_TIME|DEP_TIME|WHEELS_ON|TAXI_IN|CRS_ARR_TIME|ARR_TIME|CANCELLED|
+-----+
|2000-01-01|DL|1451|BOS|Boston, MA|ATL|Atlanta, GA|1115|1113|1343|5|1400|1348|false|
|2000-01-01|DL|1479|BOS|Boston, MA|ATL|Atlanta, GA|1315|1311|1536|7|1559|1543|false|
|2000-01-01|DL|1857|BOS|Boston, MA|ATL|Atlanta, GA|1415|1414|1642|9|1721|1651|false|
|2000-01-01|DL|1997|BOS|Boston, MA|ATL|Atlanta, GA|1715|1720|1955|10|2013|2005|false|
|2000-01-01|DL|2065|BOS|Boston, MA|ATL|Atlanta, GA|2015|2010|2230|10|2300|2240|false|
|2000-01-01|US|2619|BOS|Boston, MA|ATL|Atlanta, GA|650|649|956|7|955|1003|false|
|2000-01-01|US|2621|BOS|Boston, MA|ATL|Atlanta, GA|1440|1446|1713|4|1738|1717|false|
|2000-01-01|DL|346|BTR|Baton Rouge, LA|ATL|Atlanta, GA|1740|1744|1957|9|2008|2006|false|
|2000-01-01|DL|412|BTR|Baton Rouge, LA|ATL|Atlanta, GA|1345|1345|1552|9|1622|1601|false|
|2000-01-01|DL|299|BUF|Buffalo, NY|ATL|Atlanta, GA|1245|1245|1443|5|1455|1448|false|
+-----+
only showing top 10 rows
Press Enter to Finish...←
```

Spark UI is available to your localhost:4040 port. So let me start the browser and visit localhost:4040. You can see the same Spark UI on your local machine. This UI is only available if your Spark application is running. As soon as your application finishes, the Spark UI goes away.

The screenshot shows the Apache Spark 3.2.1 Jobs UI interface. The title bar says "WritingDataFrame - Spark Jobs". The address bar shows "localhost:4040/jobs/". The top navigation bar includes links for Jobs, Stages, Storage, Environment, Executors, and SQL. The main content area is titled "Spark Jobs" with a link to "Event Timeline". Below it, it displays user information: "User: Satyam Kumar", "Total Uptime: 2.6 min", "Scheduling Mode: FIFO", and "Completed Jobs: 1". A section titled "Completed Jobs (1)" is expanded, showing a table with one row. The table columns are Job Id, Description, Submitted, Duration, Stages: Succeeded/Total, and Tasks (for all stages): Succeeded/Total. The single row has Job Id 0, Description "showString at NativeMethodAccessorImpl.java:0", Submitted "2022/05/04 11:50:34", Duration "3 s", Stages: Succeeded/Total "1/1", and Tasks (for all stages): Succeeded/Total "1/1". Navigation controls at the bottom allow for page selection and item count per page.

Job Id	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
0	showString at NativeMethodAccessorImpl.java:0 showString at NativeMethodAccessorImpl.java:0	2022/05/04 11:50:34	3 s	1/1	1/1



Thank You  
ScholarNest Technologies Pvt Ltd.  
[www.scholarnest.com](http://www.scholarnest.com)