



ScholarNest

Spark Azure Databricks

Databricks Spark Certification and beyond

Instructor: Prashant Kumar Pandey



Absolute Beginner to Specialization in Apache Spark and Azure Databricks





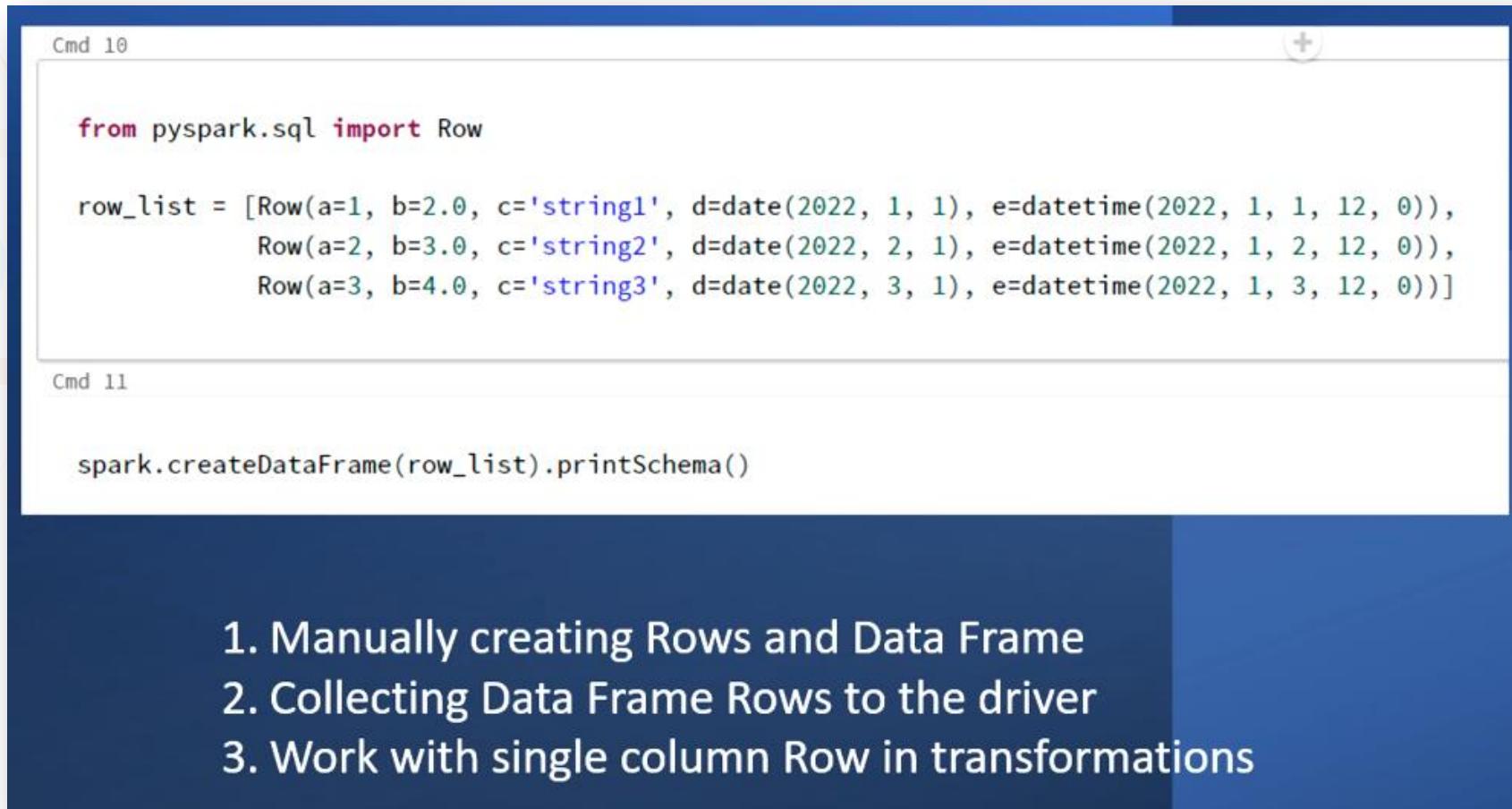
Spark Data Frame Rows

Spark DataFrame is similar to a table of rows and columns. A row in a DataFrame is a single record represented by an object of type Row. Most of the time, you do not directly work with the entire row. However, we have three specific scenarios when you might have to work with the Row object:

1. Manually creating Rows and DataFrame.
2. Collecting DataFrame rows to the driver.
3. Work with single column Rows in transformations

The first two scenarios are primarily used in unit testing or during development. And we already learned the first scenario in an earlier lecture. Here is a screenshot of the earlier example shown below.

We manually created a list of Row objects and used it to create a DataFrame.



```
from pyspark.sql import Row

row_list = [Row(a=1, b=2.0, c='string1', d=date(2022, 1, 1), e=datetime(2022, 1, 1, 12, 0)),
            Row(a=2, b=3.0, c='string2', d=date(2022, 2, 1), e=datetime(2022, 1, 2, 12, 0)),
            Row(a=3, b=4.0, c='string3', d=date(2022, 3, 1), e=datetime(2022, 1, 3, 12, 0))]

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

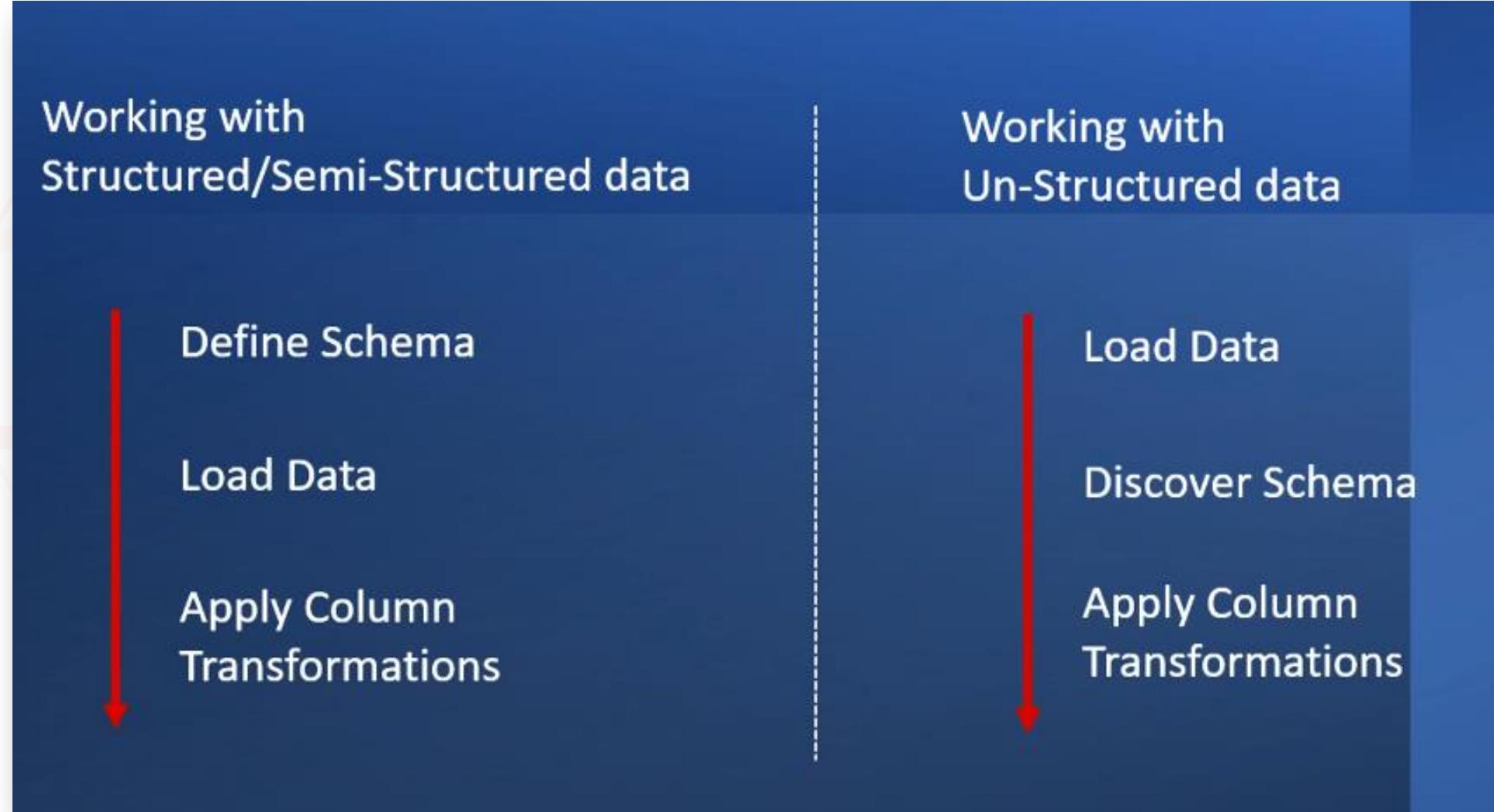
1. Manually creating Rows and Data Frame
2. Collecting Data Frame Rows to the driver
3. Work with single column Row in transformations

The third scenario of “working with single column rows in transformations” is mainly used with unstructured data.

I have this log file shown below. This file contains the application log from the Apache web server. This file contains a lot of information. But we are interested in the HTTP:// type of URLs. These are referrers of our website. They referred users to your website, and you want to take them out from this log and count who sent how many users.

```
83.149.9.216 - - [17/May/2015:10:05:03 +0000] "GET /presentations/logstash-monitorama-2013/images/kibana-search.png HTTP/1.1" 200 203023
"HTTP://semicomplete.com/presentations/logstash-monitorama-2013/" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_1) AppleWebKit/537.36 (KHTML,
Chrome/32.0.1700.77 Safari/537.36"
83.149.9.216 - - [17/May/2015:10:05:43 +0000] "GET /presentations/logstash-monitorama-2013/images/kibana-dashboard3.png HTTP/1.1" 200 171717
"HTTP://semicomplete.com/presentations/logstash-monitorama-2013/" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_1) AppleWebKit/537.36 (KHTML,
Chrome/32.0.1700.77 Safari/537.36"
83.149.9.216 - - [17/May/2015:10:05:47 +0000] "GET /presentations/logstash-monitorama-2013/plugin/highlight/highlight.js HTTP/1.1" 200 26185
"HTTP://semicomplete.com/presentations/logstash-monitorama-2013/" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_1) AppleWebKit/537.36 (KHTML,
Chrome/32.0.1700.77 Safari/537.36"
83.149.9.216 - - [17/May/2015:10:05:12 +0000] "GET /presentations/logstash-monitorama-2013/plugin/zoom-js/zoom.js HTTP/1.1" 200 7697
"HTTP://semicomplete.com/presentations/logstash-monitorama-2013/" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_1) AppleWebKit/537.36 (KHTML,
Chrome/32.0.1700.77 Safari/537.36"
83.149.9.216 - - [17/May/2015:10:05:07 +0000] "GET /presentations/logstash-monitorama-2013/plugin/notes/notes.js HTTP/1.1" 200 2892
"HTTP://semicomplete.com/presentations/logstash-monitorama-2013/" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_1) AppleWebKit/537.36 (KHTML,
Chrome/32.0.1700.77 Safari/537.36"
83.149.9.216 - - [17/May/2015:10:05:34 +0000] "GET /presentations/logstash-monitorama-2013/images/sad-medic.png HTTP/1.1" 200 430406
"HTTP://semicomplete.com/presentations/logstash-monitorama-2013/" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_1) AppleWebKit/537.36 (KHTML,
Chrome/32.0.1700.77 Safari/537.36"
83.149.9.216 - - [17/May/2015:10:05:57 +0000] "GET /presentations/logstash-monitorama-2013/css/fonts/Roboto-Bold.ttf HTTP/1.1" 200 38720
"HTTP://semicomplete.com/presentations/logstash-monitorama-2013/" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_1) AppleWebKit/537.36 (KHTML,
Chrome/32.0.1700.77 Safari/537.36"
83.149.9.216 - - [17/May/2015:10:05:50 +0000] "GET /presentations/logstash-monitorama-2013/css/fonts/Roboto-Regular.ttf HTTP/1.1" 200 41820
"HTTP://semicomplete.com/presentations/logstash-monitorama-2013/" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_1) AppleWebKit/537.36 (KHTML,
Chrome/32.0.1700.77 Safari/537.36"
83.149.9.216 - - [17/May/2015:10:05:24 +0000] "GET /presentations/logstash-monitorama-2013/images/frontend-response-codes.png HTTP/1.1" 200 5
"HTTP://semicomplete.com/presentations/logstash-monitorama-2013/" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_1) AppleWebKit/537.36 (KHTML,
Chrome/32.0.1700.77 Safari/537.36"
83.149.9.216 - - [17/May/2015:10:05:50 +0000] "GET /presentations/logstash-monitorama-2013/images/kibana-dashboard.png HTTP/1.1" 200 321631
"HTTP://semicomplete.com/presentations/logstash-monitorama-2013/" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_1) AppleWebKit/537.36 (KHTML,
Chrome/32.0.1700.77 Safari/537.36"
83.149.9.216 - - [17/May/2015:10:05:46 +0000] "GET /presentations/logstash-monitorama-2013/images/Dreamhost_logo.svg HTTP/1.1" 200 2126
"HTTP://semicomplete.com/presentations/logstash-monitorama-2013/" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_1) AppleWebKit/537.36 (KHTML,
Chrome/32.0.1700.77 Safari/537.36"
83.149.9.216 - - [17/May/2015:10:05:11 +0000] "GET /presentations/logstash-monitorama-2013/images/kibana-dashboard2.png HTTP/1.1" 200 394967
"HTTP://semicomplete.com/presentations/logstash-monitorama-2013/" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_1) AppleWebKit/537.36 (KHTML,
Chrome/32.0.1700.77 Safari/537.36"
83.149.9.216 - - [17/May/2015:10:05:19 +0000] "GET /presentations/logstash-monitorama-2013/images/apache-icon.gif HTTP/1.1" 200 8095
"HTTP://semicomplete.com/presentations/logstash-monitorama-2013/" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_1) AppleWebKit/537.36 (KHTML,
Chrome/32.0.1700.77 Safari/537.36"
```

Spark allows you to work with structured/semi-structured and unstructured data. However, the approach is different as shown below.



When you have structured or semi-structured data, you will define the schema of your data source.

Load it to a Spark Dataframe.

And start applying transformations.

But what if you have unstructured data.

The approach is different.

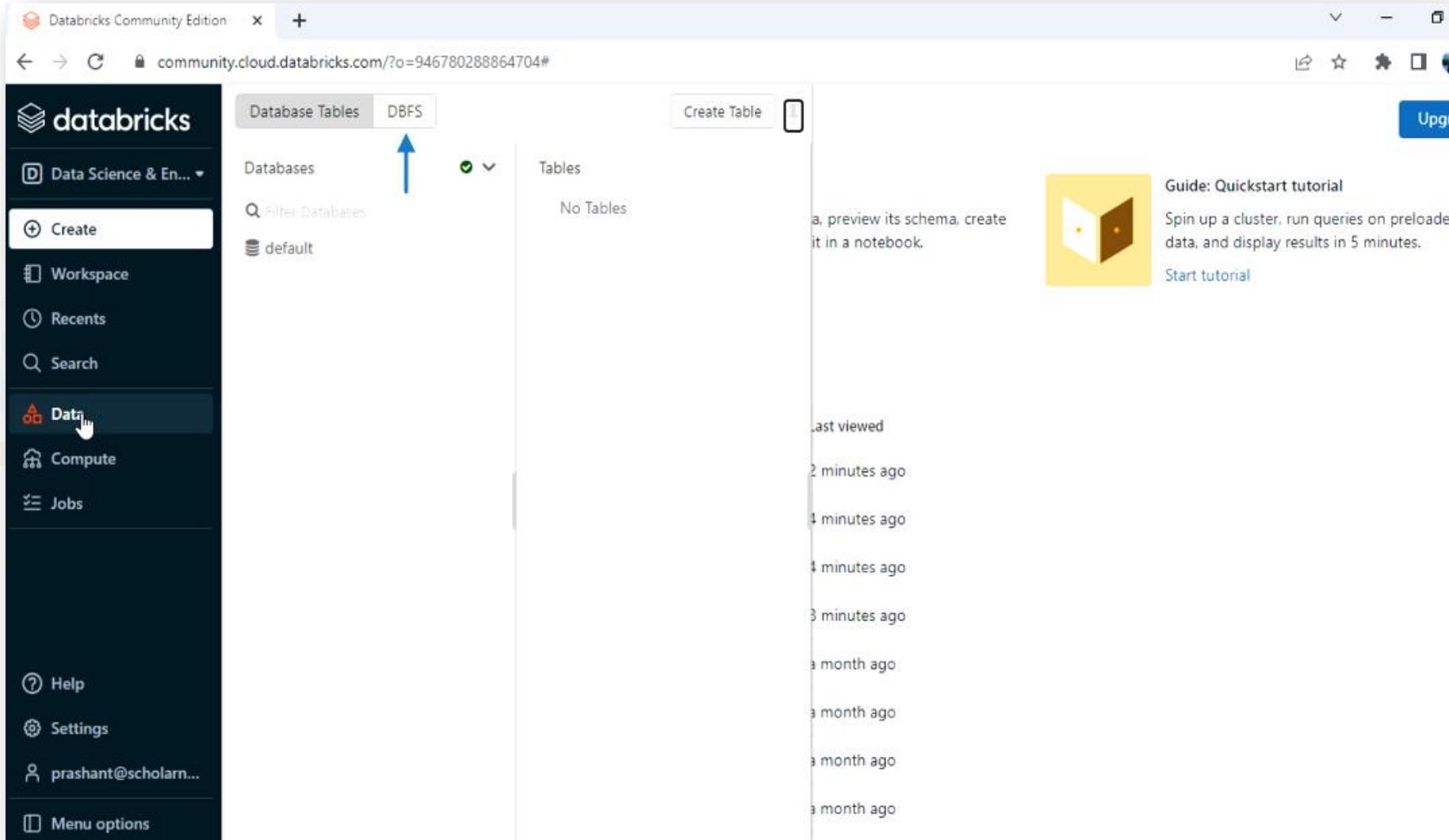
Load the unstructured data into a Dataframe. You won't have a meaningful schema.

It is most likely a line of text in each row.

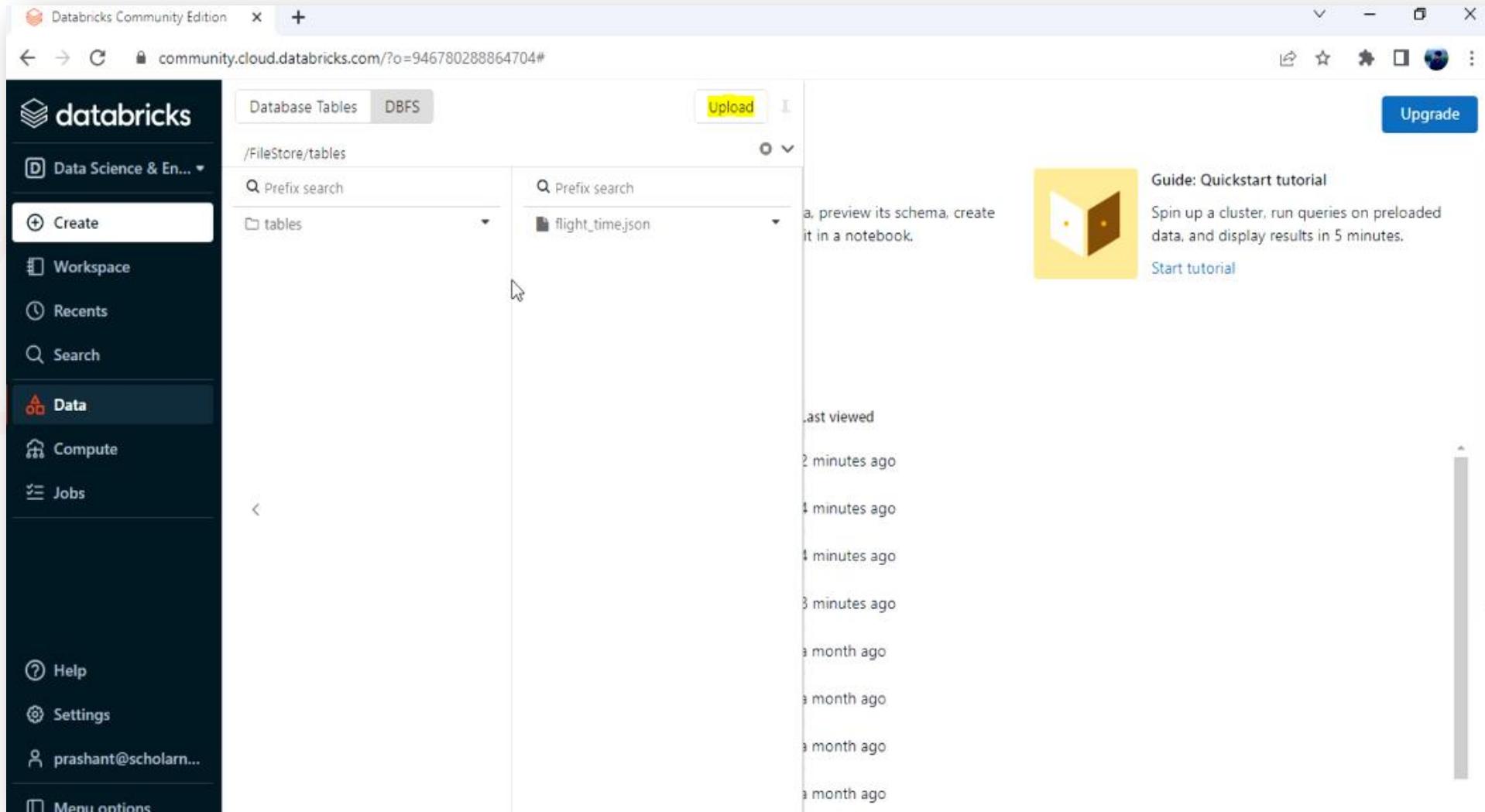
But then you can discover a schema and give your Dataframe a row/column structure.

Basically, we convert the unstructured data into structured and start transforming and analysing it.

Go to your data bricks workspace and upload the log file into DBFS.
Go to the Data menu. And you will see the DBFS tab. Click that, and you can navigate to the filesystem.



Select Filestore/tables and then click the upload button. And it will show the DBFS target directory here. You can click the upload option and choose your data file to upload.

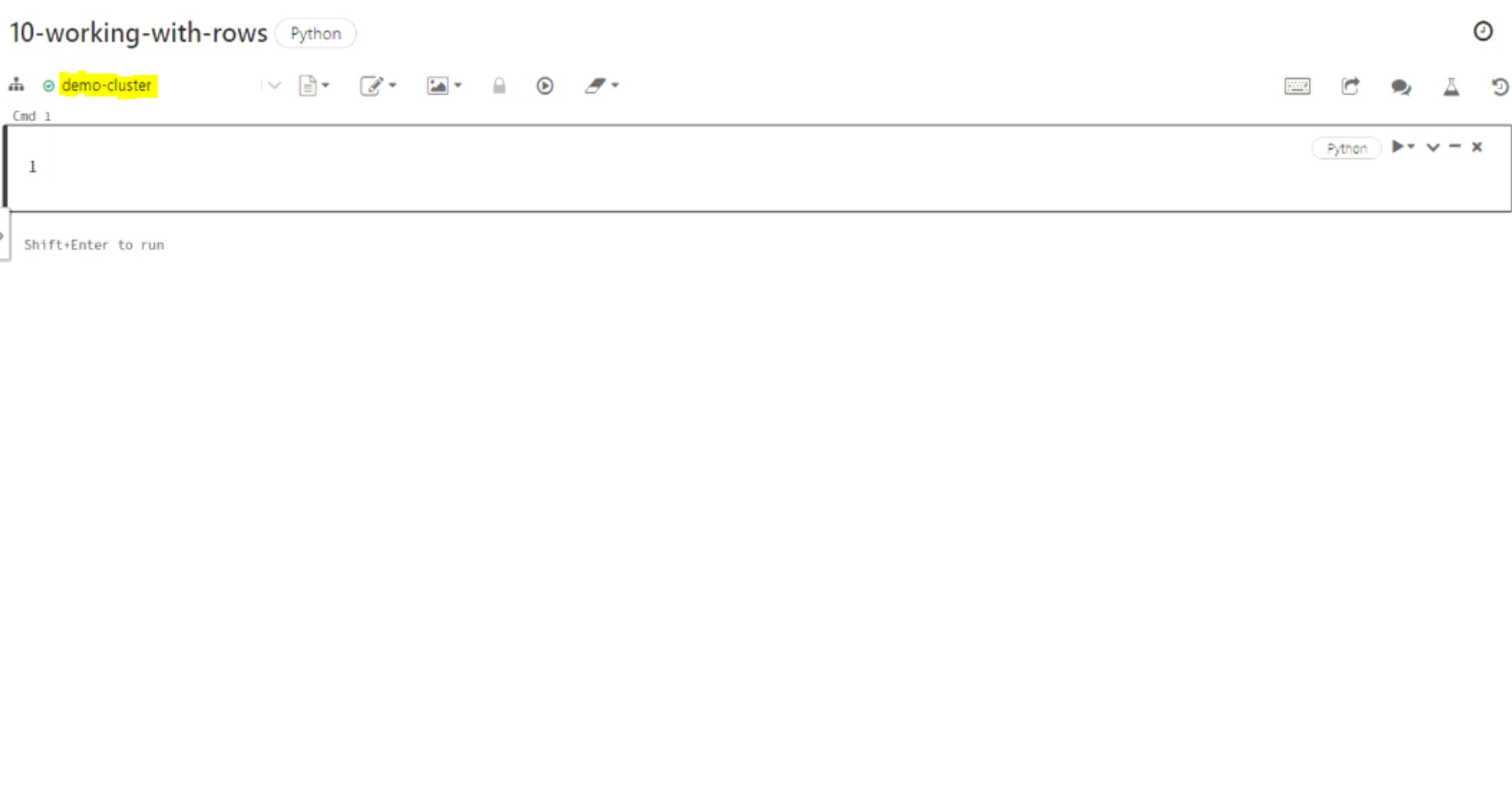


My data file is in the cloud, and we are now ready to start processing it.

The screenshot shows the Databricks Community Edition interface. On the left, there's a sidebar with options like Create, Workspace, Recents, Search, Data, Compute, and Jobs. The main area shows a 'Database Tables' tab selected, with a 'Upload' button at the top. A modal window titled 'Upload Data to DBFS' is open, showing a file named 'apache-logs.txt' uploaded to '/FileStore/tables'. The file is 2.4 MB and has a green checkmark next to it. Below the file, a message says 'File uploaded to /FileStore/tables/apache_logs.txt'. A 'Done' button is at the bottom right of the modal. In the background, there's a sidebar with a guide icon and text: 'Guide: Quickstart tutorial' and 'Spin up a cluster, run q data, and display result Start tutorial'.

Create a new notebook.(Reference : **10-working-with-rows**)

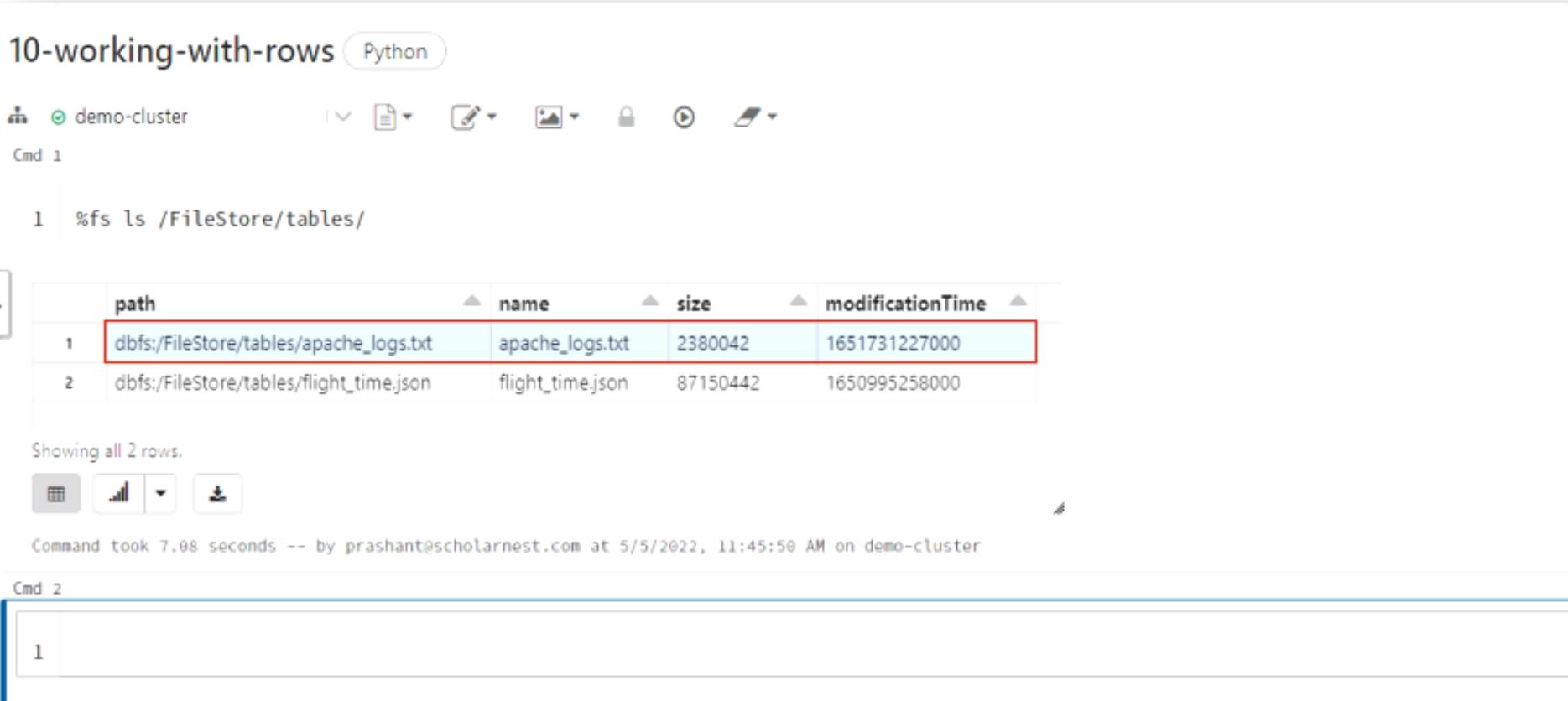
Make sure you attach a running cluster to your notebook.



The screenshot shows a Jupyter Notebook interface. The title bar says "10-working-with-rows" and "Python". Below the title bar is a toolbar with various icons. The main area has a cell labeled "Cmd 1" containing the number "1". To the left of the cell is a small icon with a right-pointing arrow. Below the cell is a prompt: "Shift+Enter to run". The background of the slide features a large, faint watermark of a red triangle.

We have already uploaded the log file to the cloud storage. You can use the command shown below to see the file.

The next step is to load this file into a Spark DataFrame.



```
10-working-with-rows Python
demo-cluster
Cmd 1
1 %fs ls /FileStore/tables/
>
1 | path           | name        | size   | modificationTime |
| dbfs:/FileStore/tables/apache_logs.txt | apache_logs.txt | 2380042 | 1651731227000 |
2 | dbfs:/FileStore/tables/flight_time.json | flight_time.json | 87150442 | 1650995258000 |
Showing all 2 rows.
Command took 7.08 seconds -- by prashant@scholarnest.com at 5/5/2022, 11:45:50 AM on demo-cluster
Cmd 2
1
Shift+Enter to run
```

Here is the code to load the file into Spark Dataframe.

We start with the spark session variable, use the read attribute to get the DataFrameReader, set the format, and load it. The format is text. Then, I am also printing the schema of the data frame.

And you can see the output for this code below the cell. Dataframe must have a schema. You cannot create a schema-less Dataframe. We are loading a text file which is an unstructured data file. However, the Dataframe must have a schema even if you load an unstructured data file.

The screenshot shows a Jupyter Notebook cell titled "10-working-with-rows" in Python. The cell contains the following code:

```
1 file_df = spark.read \
2     .format("text") \
3     .load("/FileStore/tables/apache_logs.txt")
4
5 file_df.printSchema()
```

The output of the code is:

```
root
 |-- value: string (nullable = true)
```

At the bottom of the cell, there is a note: "Command took 1.69 seconds -- by prashant@scholarnest.com at 5/5/2022, 11:47:42 AM on demo-cluster".

Here is the Dataframe. We have a single-column data frame. The column name is value. Each record comes with one line of text, and it is an unstructured Dataframe. Now we want to do some analysis on this data and apply Dataframe transformations. However, most of the Dataframe transformations are applied to columns. And we do not have any column except the value column, because we have unstructured data.

10-working-with-rows Python

demo-cluster

1 display(file_df)

(1) Spark Jobs

	value
1	83.149.9.216 - - [17/May/2015:10:05:03 +0000] "GET /presentations/logstash-monitorama-2013/images/kibana-search.png HTTP/1.1" 200 203023 "http://semicomplete.com/presentations/logstash-monitorama-2013/" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_1) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/32.0.1700.77 Safari/537.36"
2	83.149.9.216 - - [17/May/2015:10:05:43 +0000] "GET /presentations/logstash-monitorama-2013/images/kibana-dashboard3.png HTTP/1.1" 200 171717 "http://semicomplete.com/presentations/logstash-monitorama-2013/" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_1) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/32.0.1700.77 Safari/537.36"
3	83.149.9.216 - - [17/May/2015:10:05:47 +0000] "GET /presentations/logstash-monitorama-2013/plugin/highlight/highlight.js HTTP/1.1" 200 26185 "http://semicomplete.com/presentations/logstash-monitorama-2013/" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_1) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/32.0.1700.77 Safari/537.36"
4	83.149.9.216 - - [17/May/2015:10:05:12 +0000] "GET /presentations/logstash-monitorama-2013/plugin/zoom-js/zoom.js HTTP/1.1" 200 7697 "http://semicomplete.com/presentations/logstash-monitorama-2013/" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_1) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/32.0.1700.77 Safari/537.36"
5	83.149.9.216 - - [17/May/2015:10:05:07 +0000] "GET /presentations/logstash-monitorama-2013/plugin/notes/notes.js HTTP/1.1" 200 2892 "http://semicomplete.com/presentations/logstash-monitorama-2013/" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_1) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/32.0.1700.77 Safari/537.36"
6	83.149.9.216 - - [17/May/2015:10:05:34 +0000] "GET /presentations/logstash-monitorama-2013/images/sad-medic.png HTTP/1.1" 200 430406 "http://semicomplete.com/presentations/logstash-monitorama-2013/" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_1) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/32.0.1700.77 Safari/537.36"
	83.149.9.216 - - [17/May/2015:10:05:57 +0000] "GET /presentations/logstash-monitorama-2013/css/fonts/Roboto-Bold.ttf

Truncated results, showing first 1000 rows.

I gave some structure to the data, by taking one record, and breaking it with spaces. So I identified 11 fields in this data. The first one is the IP address. The second and third ones are hyphens. The fourth one is a date. Then we have Get or Post method. The next one is an absolute file name. Then we have protocol followed by two numbers. After all this, we have the referrer URL I am looking for. Once I get a URL, I do not care what else is there, so consider it one field.

```
83.149.9.216 - - [17/May/2015:10:05:03 +0000] "GET /presentations/logstash-
monitorama-2013/images/kibana-search.png HTTP/1.1" 200 203023
"http://semicomplete.com/presentations/logstash-monitorama-2013/" "Mozilla/5.0
(Macintosh; Intel Mac OS X 10_9_1) AppleWebKit/537.36 (KHTML, like Gecko)
Chrome/32.0.1700.77 Safari/537.36"
```

1. 83.149.9.216
2. -
3. -
4. [17/May/2015:10:05:03 +0000]
5. "GET
6. /presentations/logstash-monitorama-2013/images/kibana-search.png
7. HTTP/1.1"
8. 200
9. 203023
10. "http://semicomplete.com/presentations/logstash-monitorama-2013/"
11. "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_1) AppleWebKit/537.36 (KHTML,
like Gecko) Chrome/32.0.1700.77 Safari/537.36"

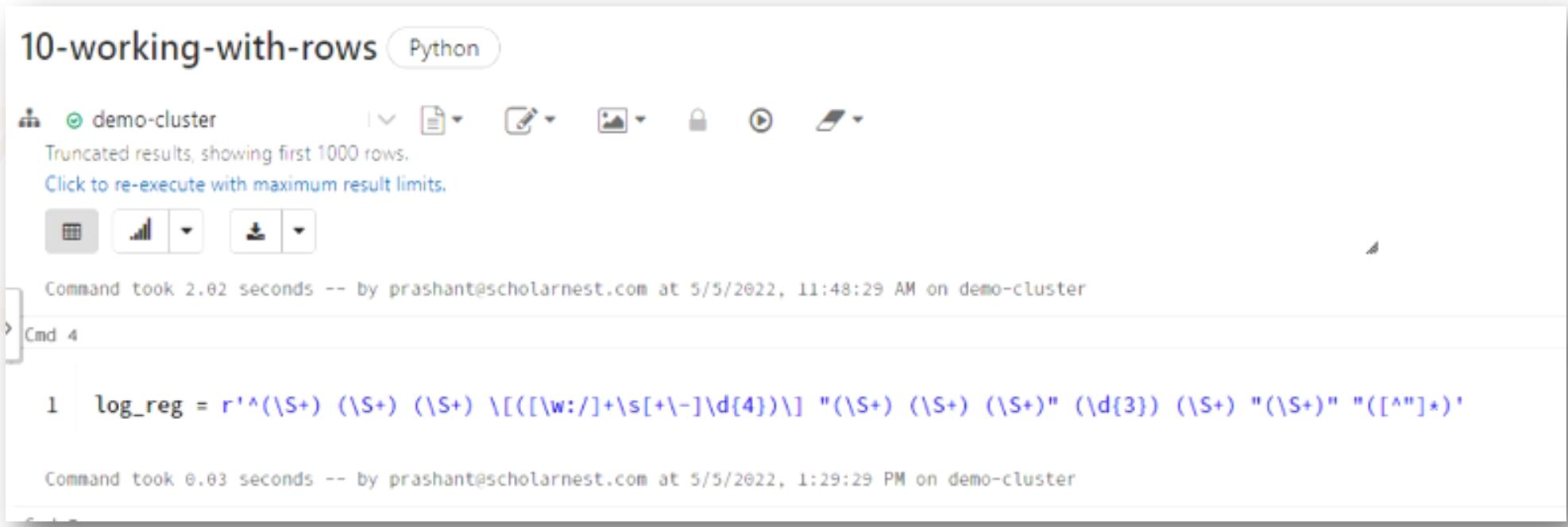
I saw the data and realized that the referrer's URL is in the file. However, we have another nine logical fields before the referrer. So If I can break this text line into eleven fields as listed here, I can put a structure to the given unstructured Dataframe. It is not easy to split this text using space characters. I thought of applying regular expressions. And here are eleven expressions assigned for each field.

```
83.149.9.216 - - [17/May/2015:10:05:03 +0000] "GET /presentations/logstash-monitorama-2013/images/kibana-search.png HTTP/1.1" 200 203023  
"http://semicomplete.com/presentations/logstash-monitorama-2013/" "Mozilla/5.0  
(Macintosh; Intel Mac OS X 10_9_1) AppleWebKit/537.36 (KHTML, like Gecko)  
Chrome/32.0.1700.77 Safari/537.36"
```

1. 83.149.9.216
2. -
3. -
4. [17/May/2015:10:05:03 +0000]
5. "GET
6. /presentations/logstash-monitorama-2013/images/kibana-search.png
7. HTTP/1.1"
8. 200
9. 203023
10. "http://semicomplete.com/presentations/logstash-monitorama-2013/"
11. "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_1) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/32.0.1700.77 Safari/537.36"

1. (\S+)
2. (\S+)
3. (\S+)
4. \[([{\w:/}]+\s[+\-]\d{4})\]
5. "(\S+)"
6. (\S+)
7. (\S+)"
8. (\d{3})
9. (\S+)
10. "(\S+)"
11. "([^\"]*)"

Here is the regular expression that represents all the eleven fields.



The screenshot shows a Jupyter Notebook cell titled "10-working-with-rows" in Python mode. The cell content is a regular expression definition:

```
1 log_reg = r'^(\S+) (\S+) (\S+) \[(\w:/]+\s[+\-]\d{4})\] "(\S+) (\S+) (\S+)" (\d{3}) (\S+) "(\S+)" "([^\"]*)"'
```

Below the code, the output shows the command took 0.03 seconds to execute on a demo-cluster at 5/5/2022, 1:29:29 PM.

Spark gives us `regexp_extract()` function to extract one field from a regular expression, and I have the code for that shown below. So I am using the `select()` method to take out four fields.

The `regexp_extract()` method takes three arguments. The first one is the column on which you want to apply the regular expression. The second argument is the regular expression itself.

The last argument is the position after splitting the value using the regular expression.

10-working-with-rows Python

demo-cluster

```
1 from pyspark.sql.functions import regexp_extract
2
3 logs_df = file_df.select(regexp_extract('value', log_reg, 1).alias('ip'),
4                           regexp_extract('value', log_reg, 4).alias('date'),
5                           regexp_extract('value', log_reg, 6).alias('image'),
6                           regexp_extract('value', log_reg, 10).alias('referrer'))
7
8 display(logs_df)
```

(1) Spark Jobs

	ip	date	image	referrer
1	83.149.9.216	17/May/2015:10:05:03 +0000	/presentations/logstash-monitorama-2013/images/kibana-search.png	http://semicom
2	83.149.9.216	17/May/2015:10:05:43 +0000	/presentations/logstash-monitorama-2013/images/kibana-dashboard3.png	http://semicom
3	83.149.9.216	17/May/2015:10:05:47 +0000	/presentations/logstash-monitorama-2013/plugin/highlight/highlight.js	http://semicom
4	83.149.9.216	17/May/2015:10:05:12 +0000	/presentations/logstash-monitorama-2013/plugin/zoom-js/zoom.js	http://semicom
5	83.149.9.216	17/May/2015:10:05:07 +0000	/presentations/logstash-monitorama-2013/plugin/notes/notes.js	http://semicom
6	83.149.9.216	17/May/2015:10:05:34 +0000	/presentations/logstash-monitorama-2013/images/sad-medic.png	http://semicom
7	83.149.9.216	17/May/2015:10:05:57 +0000	/presentations/logstash-monitorama-2013/css/fonts/Roboto-Bold.ttf	http://semicom

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

Command took 0.84 seconds -- by prashant@scholarnest.com at 5/5/2022, 1:50:56 PM on demo-cluster

I know the first position is for IP address, the fourth is a date, the sixth is the image, and the tenth field is the referrer. I am only interested in a referrer. But I tried to extract four fields so you can learn how to extract multiple files.

10-working-with-rows Python

demo-cluster Cmd 5

```
1 from pyspark.sql.functions import regexp_extract
2
3 logs_df = file_df.select(regexp_extract('value', log_reg, 1).alias('ip'),
4                           regexp_extract('value', log_reg, 4).alias('date'),
5                           regexp_extract('value', log_reg, 6).alias('image'),
6                           regexp_extract('value', log_reg, 10).alias('referrer'))
7
8 display(logs_df)
```

▶ (1) Spark Jobs

	ip	date	image	referrer
1	83.149.9.216	17/May/2015:10:05:03 +0000	/presentations/logstash-monitorama-2013/images/kibana-search.png	http://semicom
2	83.149.9.216	17/May/2015:10:05:43 +0000	/presentations/logstash-monitorama-2013/images/kibana-dashboard3.png	http://semicom
3	83.149.9.216	17/May/2015:10:05:47 +0000	/presentations/logstash-monitorama-2013/plugin/highlight/highlight.js	http://semicom
4	83.149.9.216	17/May/2015:10:05:12 +0000	/presentations/logstash-monitorama-2013/plugin/zoom-js/zoom.js	http://semicom
5	83.149.9.216	17/May/2015:10:05:07 +0000	/presentations/logstash-monitorama-2013/plugin/notes/notes.js	http://semicom
6	83.149.9.216	17/May/2015:10:05:34 +0000	/presentations/logstash-monitorama-2013/images/sad-medic.png	http://semicom
7	83.149.9.216	17/May/2015:10:05:57 +0000	/presentations/logstash-monitorama-2013/css/fonts/Roboto-Bold.ttf	http://semicom

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

Command took 0.84 seconds -- by prashant@scholarnest.com at 5/5/2022, 1:50:56 PM on demo-cluster

Here is the schema definition for the Dataframe. Now I can apply some column transformations.

10-working-with-rows Python

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

Command took 0.84 seconds -- by prashant@scholarnest.com at 5/5/2022, 1:50:56 PM on demo-cluster

> Cmd 6

```
1 logs_df.printSchema()

root
|-- ip: string (nullable = true)
|-- date: string (nullable = true)
|-- image: string (nullable = true)
|-- referrer: string (nullable = true)

Command took 0.03 seconds -- by prashant@scholarnest.com at 5/5/2022, 2:23:38 PM on demo-cluster
```

Now I have a structured Dataframe. I can easily group by the referrer, take a count, and display it.

You can see the results. And we now know who referred how many clicks.

10-working-with-rows Python

demo-cluster Cmd 7

```
1 referrer_df = logs_df.groupBy("referrer") \
2     .count()
3
4 display(referrer_df)
```

▶ (2) Spark Jobs

referrer	count
1 http://www.semicomplete.com/blog/tags/wifi	1
2 http://manpages.ubuntu.com/manpages/lucid/man1/xdotool.1.html	2
3 http://s-chassis.co.nz/viewtopic.php?f=16&t=9265&sid=fad3ad046ddaa63e65cda6e4e0033c19&start=200	1
4 https://www.google.sk/	1
5 http://www.semicomplete.com/projects/keynav/	65
6 http://superuser.com/questions/355151/is-it-possible-to-remap-the-back-forward-keys-of-a-thinkpad-usb-keyboard-on-linu	1
7 https://www.google.ca/	3

Showing all 627 rows.

Command took 1.94 seconds -- by prashant@scholarnest.com at 5/5/2022, 2:33:45 PM on demo-cluster

Cmd 8

We have two problems in this Dataframe as listed below.

1. Look at the results, and you will see the same website not being aggregated together. I mean, these two URLs are the same websites. However, they are aggregating as two different websites.

10-working-with-rows Python

demo-cluster Cmd 7

```
1 referrer_df = logs_df.groupBy("referrer") \
2     .count()
3
4 display(referrer_df)
```

→ 1. The same website not being aggregated together.
2. You will see nulls or hyphens in the referrer field.

referrer	count
http://www.semicomplete.com/blog/tags/wifi	1
http://manpages.ubuntu.com/manpages/lucid/man1/xdotool.1.html	2
http://s-chassis.co.nz/viewtopic.php?f=16&t=9265&sid=fad3ad046ddaa63e65cda6e4e0033c19&start=200	1
https://www.google.sk/	1
http://www.semicomplete.com/projects/keynav/	65
http://superuser.com/questions/355151/is-it-possible-to-remap-the-back-forward-keys-of-a-thinkpad-usb-keyboard-on-linu	1
https://www.google.ca/	3

Showing all 627 rows.

Command took 1.94 seconds -- by prashant@scholarnest.com at 5/5/2022, 2:33:45 PM on demo-cluster

2. I have another problem with this data. Scroll down, and you will see nulls or hyphens in the referrer field. Here are 4K plus hits that are not coming from any referrer.

10-working-with-rows Python

demo-cluster Cmd 7

```
1 referrer_df = logs_df.groupBy("referrer") \
2     .count()
3
4 display(referrer_df)
```

(2) Spark Jobs

referrer	count
http://www.google.co.in/url?	1
309 sa=t&rct=j&q=&esrc=s&source=web&cd=4&ved=0CDQQFjAD&url=http%3A%2F%semicomplete.com%2Ffiles%2Flogstash%2F&ei=BF4CU4OEIYaiiAf184CA8w&usg=AFQjCNEM6iuqm7sxu8XzgkypI8Z17jGTzQ&bvm=bv.61535280,d.aGc&cad=rja	1
310 -	4073
http://www.google.com/url?	1
311 sa=t&rct=j&q=&esrc=s&source=web&cd=1&ved=0CCUQFjAA&url=http%3A%2F%2Fwww.semicomplete.com%2Fprojects%2Fdotool%2F&ei=cZABU8nwApXroASpkICADA&usg=AFQjCNE3V_aCf3-gfNcbS924S6jZ6FqffA&sig2=Gkp9GH13hT3n4yfztrphjg&bvm=bv.61535280,d.cGU	1

Showing all 627 rows.

Command took 1.94 seconds -- by prashant@scholarnest.com at 5/5/2022, 2:33:45 PM on demo-cluster

Cmd 8

Can you solve these problems listed below and give us the top five referrers? I leave it as an exercise for you.

The solution is included in the notebook source code. So you can refer to it.

1. The same website not being aggregated together.
2. You will see nulls or hyphens in the referrer field.
3. Display the top five referrers.



Thank You
ScholarNest Technologies Pvt Ltd.
www.scholarnest.com



ScholarNest

Spark Azure Databricks

Databricks Spark Certification and beyond

Instructor: Prashant Kumar Pandey



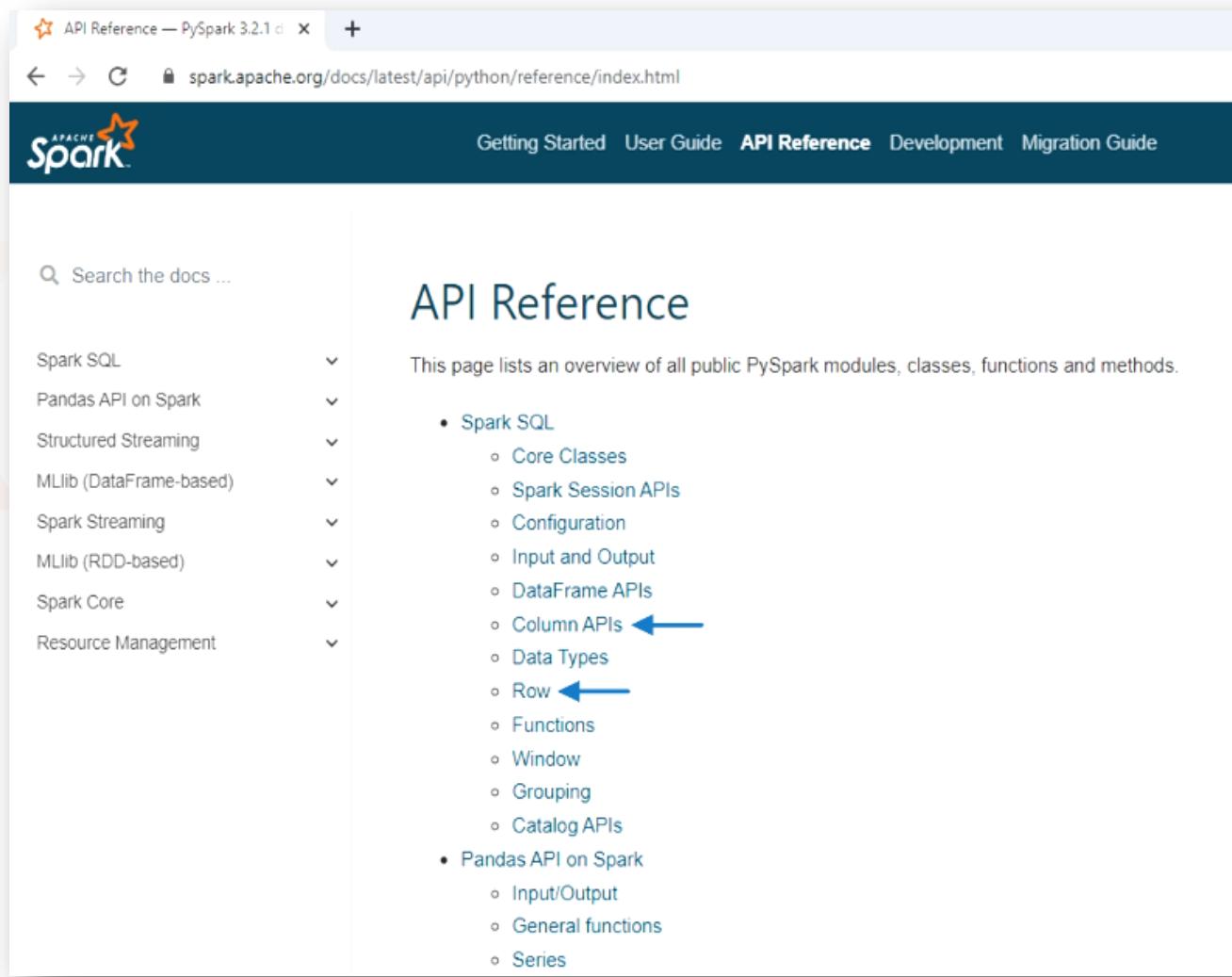
Absolute Beginner to Specialization in Apache Spark and Azure Databricks



Spark Data Frame Columns

Spark also defined two objects to represent the Dataframe Row and column. Go to the latest Spark documentation - API reference page. (**Reference -**
<https://spark.apache.org/docs/latest/api/python/reference/index.html>)

You will see two links: Column APIs and Row



The screenshot shows the Apache Spark API Reference page for PySpark 3.2.1. The page has a dark header with the Apache Spark logo and navigation links for Getting Started, User Guide, API Reference (which is bolded), Development, and Migration Guide. A search bar is at the top left. The main content area is titled "API Reference" and contains a brief overview of public PySpark modules, classes, functions, and methods. Below this, there is a list of categories under "Spark SQL" and "Pandas API on Spark". Two specific items in the "Spark SQL" list, "Column APIs" and "Row", are highlighted with blue arrows pointing to them from the left.

API Reference

This page lists an overview of all public PySpark modules, classes, functions and methods.

- [Spark SQL](#)
 - [Core Classes](#)
 - [Spark Session APIs](#)
 - [Configuration](#)
 - [Input and Output](#)
 - [DataFrame APIs](#)
 - [Column APIs](#) ←
 - [Data Types](#)
 - [Row](#) ←
 - [Functions](#)
 - [Window](#)
 - [Grouping](#)
 - [Catalog APIs](#)
- [Pandas API on Spark](#)
 - [Input/Output](#)
 - [General functions](#)
 - [Series](#)

If you go to the link for row. The row is a simple class, and it gives you only one method: `asDict()`. The `asDict()` method converts a Row into a Python dictionary. You are not likely to directly work with the Row objects. We use the Row objects in writing unit test cases for the Spark application.

The screenshot shows a web browser displaying the Apache Spark API Reference documentation. The URL in the address bar is `spark.apache.org/docs/latest/api/python/reference/pyspark.sql.html#row`. The page title is "Row". The main content area is titled "Functions" and lists several methods:

- `abs(col)`: Computes the absolute value.
- `acos(col)`: New in version 1.4.0.
- `acosh(col)`: Computes inverse hyperbolic cosine of the input column.
- `add_months(start, months)`: Returns the date that is *months* months after *start*.
- `aggregate(col, initialValue, merge[, finish])`: Applies a binary operator to an initial state and all elements in the array, and reduces this to a single state.
- `approxCountDistinct(col[, rsd])`

On the left, there is a sidebar titled "Spark SQL" with a list of related classes and functions. On the right, there is a sidebar titled "On this page" with links to other API sections like Core Classes, Spark Session APIs, Configuration, etc.

Spark Dataframe is made up of Row and Column objects. The row object is rudimentary, and we only use it in unit testing. Now let us go to the link for column API, and you will see that the column object has got many methods.

The screenshot shows a browser window displaying the Apache Spark API Reference. The title bar reads "Spark SQL — PySpark 3.2.1 docs". The address bar shows the URL "spark.apache.org/docs/latest/api/python/reference/pyspark.sql.html#column-apis". The page header includes links for "Getting Started", "User Guide", "API Reference" (which is highlighted), "Development", and "Migration Guide". The main content area is titled "Column APIs". On the left, there is a sidebar titled "Spark SQL" with a list of classes: "pyspark.sql.SparkSession", "pyspark.sql.Catalog", "pyspark.sql.DataFrame", "pyspark.sql.Column", "pyspark.sql.Row", "pyspark.sql.GroupedData", "pyspark.sql.PandasCogroupedOps", "pyspark.sql.DataFrameNaFunction", "pyspark.sql.DataFrameStatFunction", "pyspark.sql.Window", "pyspark.sql.SparkSession.builder", "pyspark.sql.SparkSession.builder.c", "pyspark.sql.SparkSession.builder.e", and "pyspark.sql.SparkSession.builder.g". The main content area lists several methods for the "Column" class:

Method	Description
<code>Column.alias(*alias, **kwargs)</code>	Returns this column aliased with a new name or names (in the case of expressions that return more than one column, such as <code>explode</code>).
<code>Column.asc()</code>	Returns a sort expression based on ascending order of the column.
<code>Column.asc_nulls_first()</code>	Returns a sort expression based on ascending order of the column, and null values return before non-null values.
<code>Column.asc_nulls_last()</code>	Returns a sort expression based on ascending order of the column, and null values appear after non-null values.
<code>Column.astype(dataType)</code>	<code>astype()</code> is an alias for <code>cast()</code> .
<code>Column.between(lowerBound, upperBound)</code>	True if the current column is between the lower bound and upper bound, inclusive.
<code>Column.bitwiseAND(other)</code>	Compute bitwise AND of this expression with another expression.

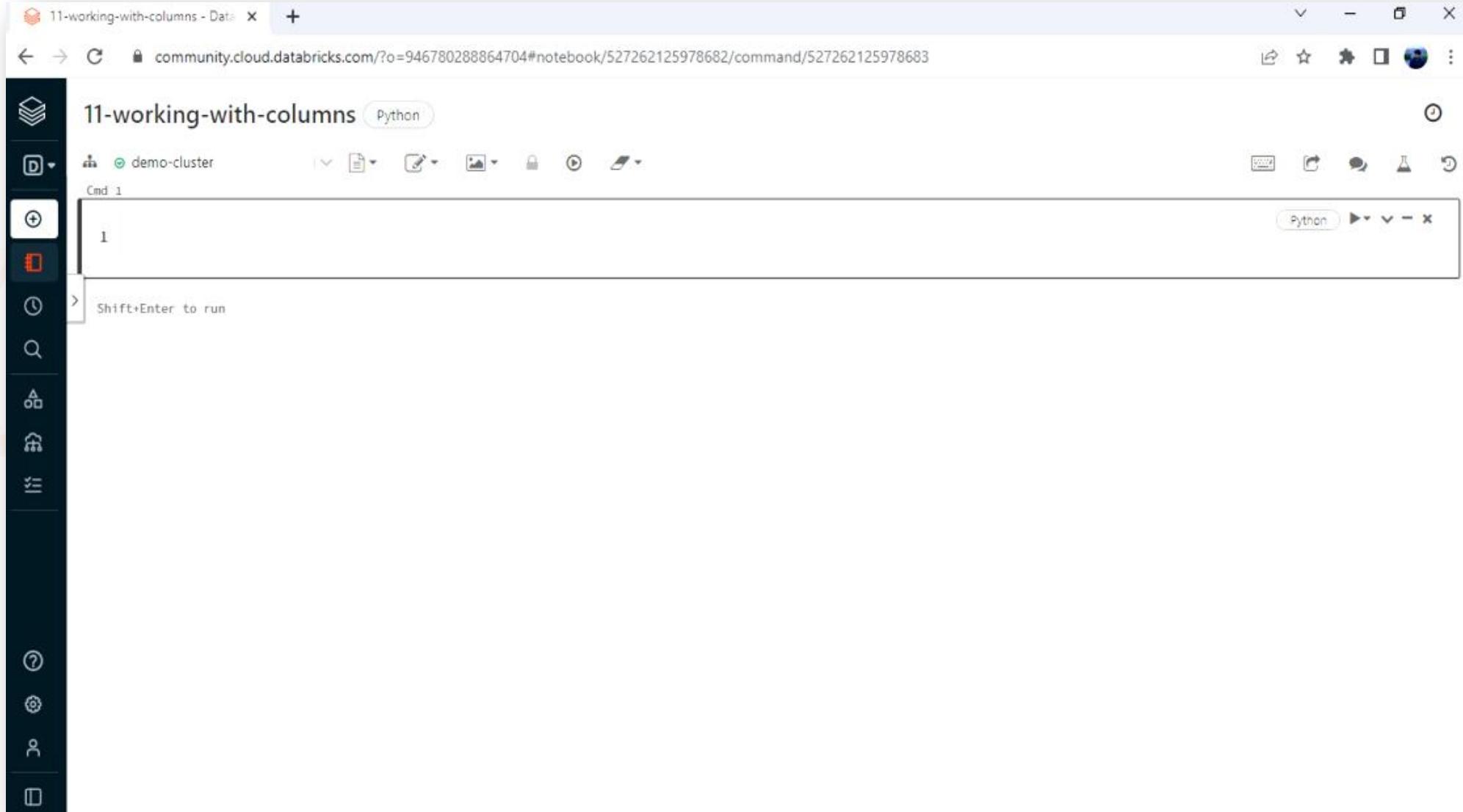
On the right side, there is a sidebar titled "On this page" with a list of links: "Core Classes", "Spark Session APIs", "Configuration", "Input and Output", "DataFrame APIs", "Column APIs" (which is highlighted), "Data Types", "Row", "Functions", "Window", "Grouping", and "Catalog APIs".

Spark Dataframe is made up of Rows and Column objects. We do not have any row transformation in spark.

And most of the Dataframe transformations are all about transforming the columns.

		Column				
		Timestamp	Age	Gender	Country	state
Row	2014-08-27 11:29:31	37	Female	United States	IL	
	2014-08-27 11:29:37	44	M	United States	IN	
	2014-08-27 11:29:44	32	Male	Canada		
	2014-08-27 11:29:46	31	Male	United Kingdom		
	2014-08-27 11:30:22	31	Male	United States	TX	

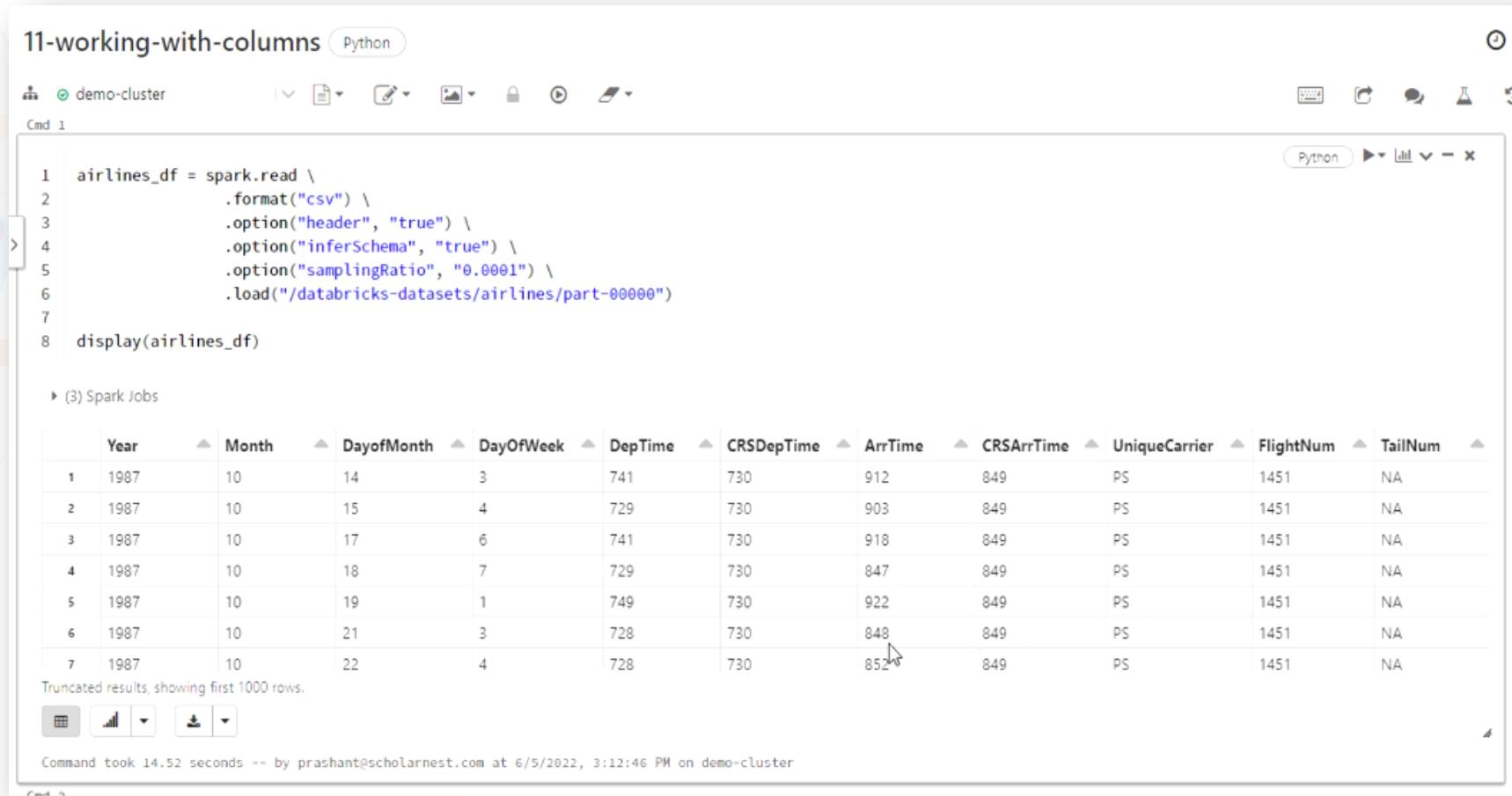
Log in to your Databricks community cloud and create a new notebook.(Reference : 11-working-with-columns)



I have created a Dataframe.

I am reading data using the spark.read. The format is CSV, and I am setting three options: Header, Infer Schema, and Sampling Ratio. Finally, I am loading the data into a Dataframe.

Infer schema option will read a sample of the data file and try to infer the schema automatically. The default value of the sample size is one percent. You can increase or decrease the sampling ratio. I am setting it to a smaller value to ensure that this command runs quickly.



The screenshot shows a Databricks notebook titled "11-working-with-columns" in Python. The code cell contains the following Python code:

```
1 airlines_df = spark.read \
2     .format("csv") \
3     .option("header", "true") \
4     .option("inferSchema", "true") \
5     .option("samplingRatio", "0.0001") \
6     .load("/databricks-datasets/airlines/part-00000")
7
8 display(airlines_df)
```

Below the code, a section titled "(3) Spark Jobs" displays a preview of the DataFrame "airlines_df". The preview shows the first 1000 rows of the dataset, which consists of 14 columns: Year, Month, DayofMonth, DayOfWeek, DepTime, CRSDepTime, ArrTime, CRSArrTime, UniqueCarrier, FlightNum, TailNum, and two additional columns that are currently NA. The data is for October 1987, with various flight details and carrier information.

	Year	Month	DayofMonth	DayOfWeek	DepTime	CRSDepTime	ArrTime	CRSArrTime	UniqueCarrier	FlightNum	TailNum	
1	1987	10	14	3	741	730	912	849	PS	1451	NA	
2	1987	10	15	4	729	730	903	849	PS	1451	NA	
3	1987	10	17	6	741	730	918	849	PS	1451	NA	
4	1987	10	18	7	729	730	847	849	PS	1451	NA	
5	1987	10	19	1	749	730	922	849	PS	1451	NA	
6	1987	10	21	3	728	730	848	849	PS	1451	NA	
7	1987	10	22	4	728	730	852	849	PS	1451	NA	

Truncated results, showing first 1000 rows.

Command took 14.52 seconds -- by prashant@scholarnest.com at 6/5/2022, 3:12:46 PM on demo-cluster

So we have an origin, destination, and distance columns.

We also have Arrival delays and departure delays before the origin column.

I want to select only these five columns.

11-working-with-columns Python

demo-cluster

Cmd 1

```
1 airlines_df = spark.read \
2     .format("csv") \
3     .option("header", "true") \
4     .option("inferSchema", "true") \
5     .option("samplingRatio", "0.0001") \
6     .load("/databricks-datasets/airlines/part-00000")
7
8 display(airlines_df)
```

(3) Spark Jobs

	ArrDelay	DepDelay	Origin	Dest	Distance	TaxiIn	TaxiOut	Cancelled	CancellationCode	Divert
1	23	11	SAN	SFO	447	NA	NA	0	NA	0
2	14	-1	SAN	SFO	447	NA	NA	0	NA	0
3	29	11	SAN	SFO	447	NA	NA	0	NA	0
4	-2	-1	SAN	SFO	447	NA	NA	0	NA	0
5	33	19	SAN	SFO	447	NA	NA	0	NA	0
6	-1	-2	SAN	SFO	447	NA	NA	0	NA	0

Truncated results, showing first 1000 rows.

Command took 14.52 seconds -- by prashant@scholarnest.com at 6/5/2022, 3:12:46 PM on demo-cluster

Here is the code to select the desired 5 columns. And we can see the output for the same given below.

11-working-with-columns Python

demo-cluster

Cmd 2

```
1 airlines_df.select("Origin", "Dest", "Distance", "ArrDelay", "DepDelay") \
2     .show(5)
```

▶ (1) Spark Jobs

Origin	Dest	Distance	ArrDelay	DepDelay
SAN	SFO	447	23	11
SAN	SFO	447	14	-1
SAN	SFO	447	29	11
SAN	SFO	447	-2	-1
SAN	SFO	447	33	19

only showing top 5 rows

Command took 0.98 seconds -- by prashant@scholarnest.com at 6/5/2022, 3:15:59 PM on demo-cluster

We have five approaches to referring to a column name.

1. "Origin"
2. `expr("Origin")`
3. `airlines_df.Origin`
4. `column("Origin")`
5. `col("Origin")`

The first approach is a column name string. The second approach uses the `expr()` function and the column name. The third one is using a `dataframe.column` name. The fourth approach is using a `column()` function. The last one is using the `col()` function.

The fourth and the fifth approaches are the same because the `col()` function is an alias for the `column()` function. So the fourth and the fifth approaches are the same.

But officially, we have five different approaches to refer to a column name. You can use any of these approaches in your select expression.

Here is the code to select the desired 5 columns. But I am referring to each column differently. I am referring to the Origin column using a plain string name. But I am referring to the destination using the expr() function. Similarly, I am referring to the distance using the dataframe.column name. The next one is using the column() function. And the last one is using the col() function.

11-working-with-columns Python

demo-cluster

Cmd 3

```
1 from pyspark.sql.functions import expr, column, col
2
3 airlines_df.select("Origin", expr("Dest"), airlines_df.Distance, column("ArrDelay"), col("DepDelay")) \
4     .show(5)
```

(1) Spark Jobs

Origin	Dest	Distance	ArrDelay	DepDelay
SAN	SFO	447	23	11
SAN	SFO	447	14	-1
SAN	SFO	447	29	11
SAN	SFO	447	-2	-1
SAN	SFO	447	33	19

only showing top 5 rows

Command took 0.60 seconds -- by prashant@scholarnest.com at 6/5/2022, 3:38:13 PM on demo-cluster

1. "Origin"
2. expr("Origin")
3. airlines_df.Origin
4. column("Origin")
5. col("Origin")

The first approach is for simplicity. I prefer using the first approach and simply writing the column name string. It is easy to use and looks neat.

But the first approach doesn't allow writing expressions.
And that's why we have the second approach of using the `expr()` function.

I am selecting three columns here in the code shown below. The first two columns are Origin and Destination. These are simple names, so they work. But the third column is a SQL expression. I am multiplying distance by 1.6 to convert it from miles to kilometers. So the third column is not a simple column name, it is an expression. You will see an error if you run it. And that's why we have the `expr()` function.

The screenshot shows a Jupyter Notebook interface with a title bar "11-working-with-columns" and a Python tab selected. The toolbar includes icons for file operations, cell execution, and help. A message bar at the top right says "Error running command 4. Go to command". Below the toolbar, a message indicates "Command took 0.60 seconds -- by prashant@scholarnest.com at 6/5/2022, 3:38:13 PM on demo-cluster". The code cell (Cmd 4) contains the following Python code:

```
> 1 airlines_df.select("Origin", "Dest", "Distance * 1.6 as km_distance") \
  2 .show(5)
```

An error message follows:

```
AnalysisException: Column ``Distance * 1.6 as km_distance`` does not exist. Did you mean one of the following? [Distance, CancellationCode, Cancelled, Diverted, UniqueCarrier, ActualElapsedTime, CRSElapsedTime, DepTime, Dest, IsArrDelayed, IsDepDelayed, LateAircraftDelay, AirTime, CRSDepTime, DayOfWeek, DayOfMonth, DepDelay, Origin, TaxiIn, WeatherDelay, ArrDelay, ArrTime, CRSArrTime, CarrierDelay, FlightNum, Month, NASDelay, SecurityDelay, TailNum, TaxiOut, Year];
'Project [Origin#89, Dest#90, 'Distance * 1.6 as km_distance']
+- Relation [Year#73,Month#74,DayofMonth#75,DayOfWeek#76,DepTime#77,CRSDepTime#78,ArrTime#79,CRSArrTime#80,UniqueCarrier#81,FlightNum#82,TailNum#83,ActualElapsedTime#84,CRSElapsedTime#85,AirTime#86,ArrDelay#87,DepDelay#88,Origin#89,Dest#90,Distance#91,TaxiIn#92,TaxiOut#93,Cancelled#94,CancellationCode#95,Diverted#96,... 7 more fields] csv
```

At the bottom, another message indicates "Command took 0.50 seconds -- by prashant@scholarnest.com at 6/5/2022, 3:39:49 PM on demo-cluster".

It works with expr() function.



11-working-with-columns Python

demo-cluster

Cmd 4

```
1 airlines_df.select("Origin", "Dest", expr("Distance * 1.6 as km_distance")) \
2     .show(5)
```

▶ (1) Spark Jobs

Origin	Dest	km_distance
SAN	SFO	715.2

only showing top 5 rows

Command took 0.83 seconds -- by prashant@scholarnest.com at 6/5/2022, 3:40:20 PM on demo-cluster

So we learned the need for the first two approaches.

- | | |
|---|---|
| 1. "Origin" | -> Used for plain column names and simplicity |
| 2. <code>expr("Distance * 1.6 as km_distance")</code> | -> Used to evaluate SQL like expression strings |

The first one is a simple column name, and we use it for simplicity.

But when you have an expression, you must evaluate it using the `expr()` function.

So the second approach is for evaluating SQL-like expression strings.

The `expr()` function also works fine for simple column names. But we do not use it for column names.

We use the `expr()` function when we have an expression.

We have the following six fields in this given Dataframe as shown in the right hand side of the screenshot. I want to write a select expression to produce the following output shown in the left side of the screenshot.

- 1. UniqueCarrier
- 2. FlightNum
- 3. Year
- 4. Month
- 5. DayofMonth
- 6. Distance

UniqueCarrier	FlightNum	FlightDate	DistanceKM
PS	1451	1987-10-14	715.2
PS	1451	1987-10-15	715.2
PS	1451	1987-10-17	715.2
PS	1451	1987-10-18	715.2
PS	1451	1987-10-19	715.2

Here is one possible solution. The first and the second columns are simply the names. So I am not using the `expr()` function for those column names. But the third and the fourth columns are expressions. So I am using the `expr()` function.

11-working-with-columns Python

demo-cluster

Cmd 5

```
1 airlines_df.select("Origin",
2                         "FlightNum",
3                         expr("Year||'-'||Month||'-'||DayOfMonth as FlightDate"),
4                         expr("Distance * 1.6 as DistanceKM")) \
5 .show(5)
```

▶ (1) Spark Jobs

Origin	FlightNum	FlightDate	DistanceKM
SAN	1451	1987-10-14	715.2
SAN	1451	1987-10-15	715.2
SAN	1451	1987-10-17	715.2
SAN	1451	1987-10-18	715.2
SAN	1451	1987-10-19	715.2

only showing top 5 rows

Command took 0.83 seconds -- by prashant@scholarnest.com at 6/5/2022, 3:43:55 PM on demo-cluster

The third approach is using the `dataframe.column` name.

We use the third approach to avoid column name ambiguity. We have another Dataframe `airport_df` shown below. The `airlines_df` is the main dataframe, and the `airport_df` is a lookup data frame. And I want to produce results like the one shown in the screenshot below.

11-working-with-columns Python

demo-cluster

Command took 0.83 seconds -- by prashant@scholarnest.com at 6/5/2022, 3:43:55 PM on demo-cluster

Cmd 6

```
1 data_list = [("SAN", "San Antonio"),
 2             ("SFO", "San Francisco"),
 3             ("BUR", "Hollywood Burbank"),
 4             ("OAK", "Oakland")]
 5
 6 airport_df = spark.createDataFrame(data_list).toDF("Origin", "OriginAirportName")
 7 display(airport_df)
```

(3) Spark Jobs

	Origin	OriginAirportName
1	SAN	San Antonio
2	SFO	San Francisco
3	BUR	Hollywood Burbank
4	OAK	Oakland

Showing all 4 rows.

Origin|OriginAirportName|Dest|Distance

Origin	OriginAirportName	Dest	Distance
BUR	Hollywood Burbank	OAK	325
BUR	Hollywood Burbank	SFO	326
BUR	Hollywood Burbank	LAS	223
BUR	Hollywood Burbank	SJC	296
BUR	Hollywood Burbank	SMF	358

Command took 1.50 seconds -- by prashant@scholarnest.com at 6/5/2022, 3:52:47 PM on demo-cluster

Cmd 7

We can use the following code to achieve the desired result. We are doing a left join of `airlines_df` with the `airport_df`. Look at the join condition highlighted below.

We have a small problem with the `airlines_df` and the `airport_df`. They both have the same column name. The column name is the origin of both the data frames. And we cannot say `Origin==Origin`, Spark will complain about ambiguous column names. So we could use the `dataframe.column name` to avoid column name ambiguity.

11-working-with-columns Python

demo-cluster

Command took 1.50 seconds -- by prashant@scholarnest.com at 6/5/2022, 3:52:47 PM on demo-cluster

Cmd 7

```
> 1  airlines_df.join(airport_df, airlines_df.Origin==airport_df.Origin, "left") \
2      .select(airlines_df.Origin, "OriginAirportName", "Dest", "Distance") \
3      .where("Origin in ('SAN', 'SFO', 'BUR', 'OAK')") \
4      .distinct() \
5      .show(5)

▶ (3) Spark Jobs
```

Origin	OriginAirportName	Dest	Distance
BUR	Hollywood Burbank	OAK	325
BUR	Hollywood Burbank	SFO	326
BUR	Hollywood Burbank	LAS	223
BUR	Hollywood Burbank	SJC	296
BUR	Hollywood Burbank	SMF	358

only showing top 5 rows

Command took 9.23 seconds -- by prashant@scholarnest.com at 6/5/2022, 3:57:36 PM on demo-cluster

Here is an expression that we created earlier. I have a cFactor variable. And I am using the variable inside the expression using Python string substitution technique. But instead of using this string technique, some people prefer to use the col() function.

11-working-with-columns Python

demo-cluster

Command took 9.23 seconds -- by prashant@scholarnest.com at 6/5/2022, 3:57:36 PM on demo-cluster

Cmd 8

```
1 cFactor = 1.6
2 airlines_df.select("Origin", "Dest", expr("Distance * {} as km_distance".format(cFactor))) \
3     .show(5)
```

▶ (1) Spark Jobs

Origin	Dest	km_distance
SAN	SFO	715.21

only showing top 5 rows

Command took 0.65 seconds -- by prashant@scholarnest.com at 6/5/2022, 4:01:16 PM on demo-cluster

Here is the same code using the col() function.

11-working-with-columns Python

demo-cluster

Cmd 9

```
1 cFactor = 1.6
2 airlines_df.select("Origin", "Dest", col("Distance") * cFactor) \
3     .show(5)
```

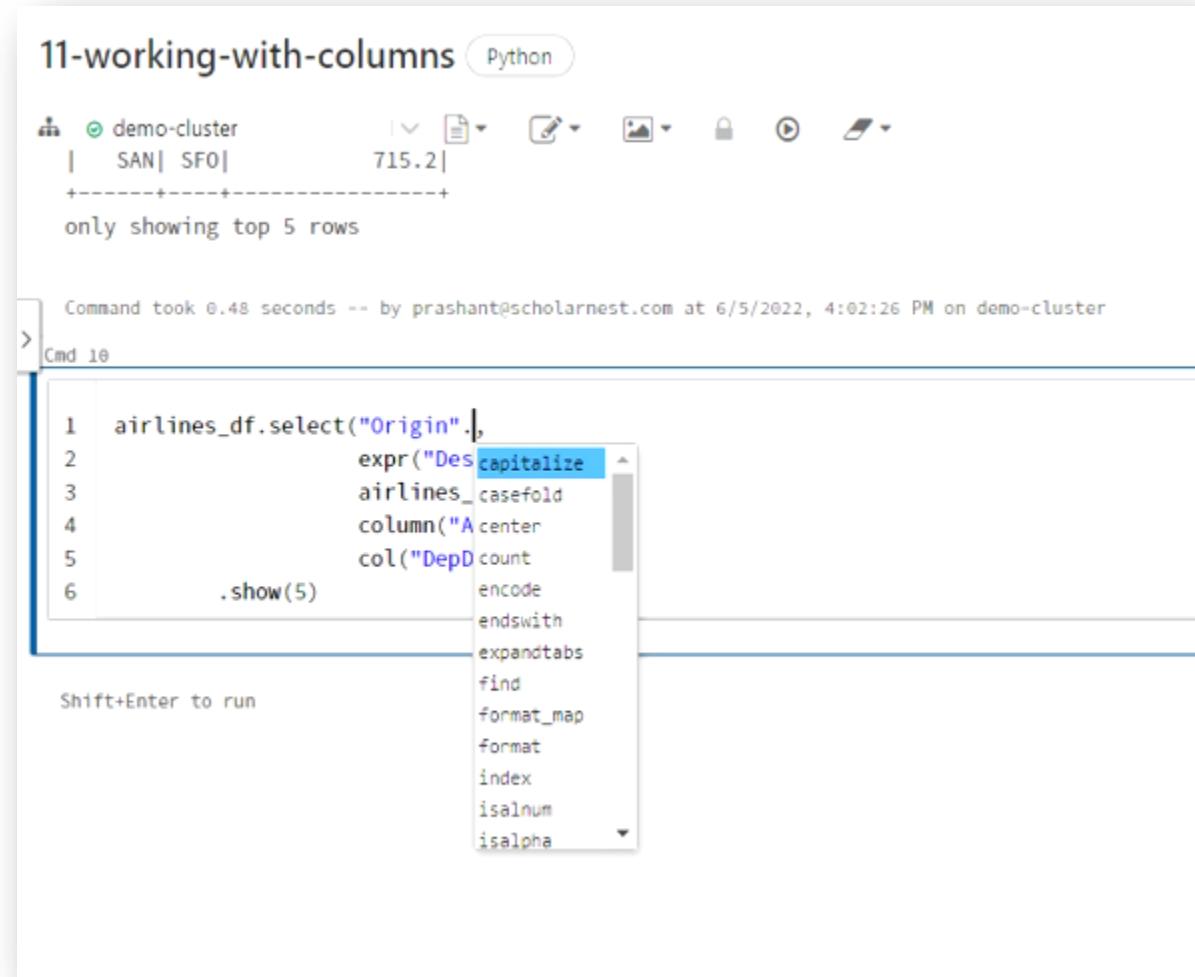
▶ (1) Spark Jobs

Origin	Dest	(Distance * 1.6)
SAN	SFO	715.2

only showing top 5 rows

Command took 0.48 seconds -- by prashant@scholarnest.com at 6/5/2022, 4:02:26 PM on demo-cluster

Here is a Dataframe select statement. I am selecting five columns and using five different approaches to select the columns. The first one is a column name string. So place your cursor after the string name, type a dot symbol and press the tab. You will see a bunch of string functions such as capitalize, casefold, center, etc. These are all Python string manipulation functions.



The screenshot shows a Jupyter Notebook cell titled "11-working-with-columns" in Python mode. The cell contains the following code:

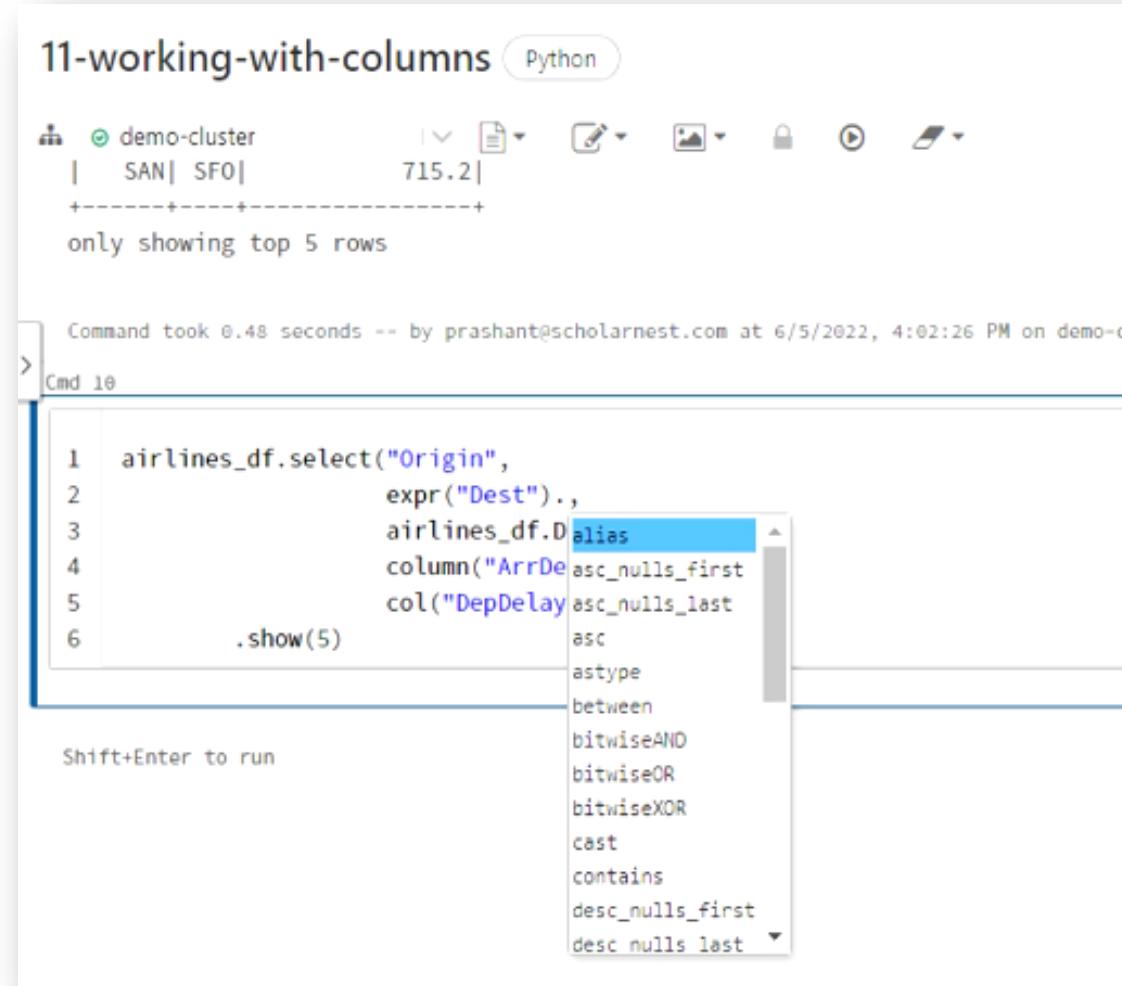
```
1 airlines_df.select("Origin").  
2     expr("Des  
3         airlines_.casefold  
4             column("A".center  
5                 col("Dep".count  
6                     .show(5)
```

A tooltip is displayed over the ".expr("Des" part of the code, listing various Python string methods starting with "expr":

- capitalise
- airlines_.casefold
- column("A".center
- col("Dep".count
- encode
- endswith
- expandtabs
- find
- format_map
- format
- index
- isalnum
- isalpha

Below the code, there is a note: "Shift+Enter to run".

But if you try getting a similar lookup for other approaches, you will see column APIs. Place your cursor at the end of the expr() function, type a dot, and press the tab. You will see column APIs such as alias, asc, between, etc. These are column APIs, and they help you create complex expressions.



The screenshot shows a Jupyter Notebook cell titled "11-working-with-columns" in Python mode. The cell contains the following code:

```
1 airlines_df.select("Origin",
2                     expr("Dest"),
3                     airlines_df.D
4                     alias
5                     column("ArrDe
6                     col("DepDelay"
7                     .show(5)
```

A code completion dropdown menu is open over the word "alias". The visible suggestions include: alias, column, desc_nulls_first, desc nulls last, asc, asc_nulls_first, asc_nulls_last, astype, between, bitwiseAND, bitwiseOR, bitwiseXOR, cast, contains, and desc nulls first. The suggestion "alias" is highlighted in blue, matching the color of the word in the code.

So I am selecting five columns here in this example. Then I am applying a where() condition. I wanted to apply two conditions.

1. where the origin is like B%
2. ArrDelay < 0

I defined both the conditions here. But I implement it using the Dataframe Column APIs. The col("Origin") will give me a Dataframe column object. The like() is a Column API. Similarly, I defined the second condition. This approach to defining expressions might look complicated for many developers. So we often do it using string expressions.

11-working-with-columns Python

demo-cluster Command took 0.48 seconds -- by prashant@scholarnest.com at 6/5/2022, 4:02:26 PM on demo-cluster

Cmd 10

```
1 airlines_df.select("Origin", "Dest", "Distance", "ArrDelay", "DepDelay") \
2     .where((col("Origin").like("B%")) &
3            (col("ArrDelay") < 0)) \
4     .show(5)
```

▶ (1) Spark Jobs

Origin	Dest	Distance	ArrDelay	DepDelay
BUR	OAK	325	-6	-1
BUR	OAK	325	-1	0
BUR	OAK	325	-6	0
BUR	OAK	325	-2	-1
BUR	OAK	325	-3	0

only showing top 5 rows

Command took 1.08 seconds -- by prashant@scholarnest.com at 6/7/2022, 2:06:02 AM on demo-cluster

Here is the same example but using the string expression. Look at the string expression here. This approach looks easy for SQL developers. The earlier approach looks easy for Python or Java Experts. But you have both options. Feel free to use any of the approaches, and Spark performs the same for both approaches. I am using the expr() function inside the where() method. But the expr() is options inside the where() method.

11-working-with-columns Python

demo-cluster only showing top 5 rows

Command took 1.08 seconds -- by prashant@scholarnest.com at 6/7/2022, 2:06:02 AM on demo-cluster

Cmd 11

```
> 1 airlines_df.select("Origin", "Dest", "Distance", "ArrDelay", "DepDelay") \
  2     .where(expr("Origin like 'B%' and ArrDelay < 0")) ←
  3     .show(5)
```

▶ (1) Spark Jobs

Origin	Dest	Distance	ArrDelay	DepDelay
BUR	OAK	325	-6	-1
BUR	OAK	325	-1	0
BUR	OAK	325	-6	0
BUR	OAK	325	-2	-1
BUR	OAK	325	-3	0

only showing top 5 rows

Command took 0.78 seconds -- by prashant@scholarnest.com at 6/7/2022, 2:07:45 AM on demo-cluster

The `where()` method expects an expression, and it will automatically evaluate the string. So you can remove the `expr()` function and simply use the string.



11-working-with-columns Python

demo-cluster only showing top 5 rows

Command took 1.08 seconds -- by prashant@scholarnest.com at 6/7/2022, 2:06:02 AM on demo-cluster

Cmd 11

```
> 1 airlines_df.select("Origin", "Dest", "Distance", "ArrDelay", "DepDelay") \
  2     .where("Origin like 'B%' and ArrDelay < 0") \
  3     .show(5)
```

▶ (1) Spark Jobs

Origin	Dest	Distance	ArrDelay	DepDelay
BUR	OAK	325	-6	-1
BUR	OAK	325	-1	0
BUR	OAK	325	-6	0
BUR	OAK	325	-2	-1
BUR	OAK	325	-3	0

only showing top 5 rows

Command took 0.86 seconds -- by prashant@scholarnest.com at 6/7/2022, 2:08:25 AM on demo-cluster



Thank You
ScholarNest Technologies Pvt Ltd.
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



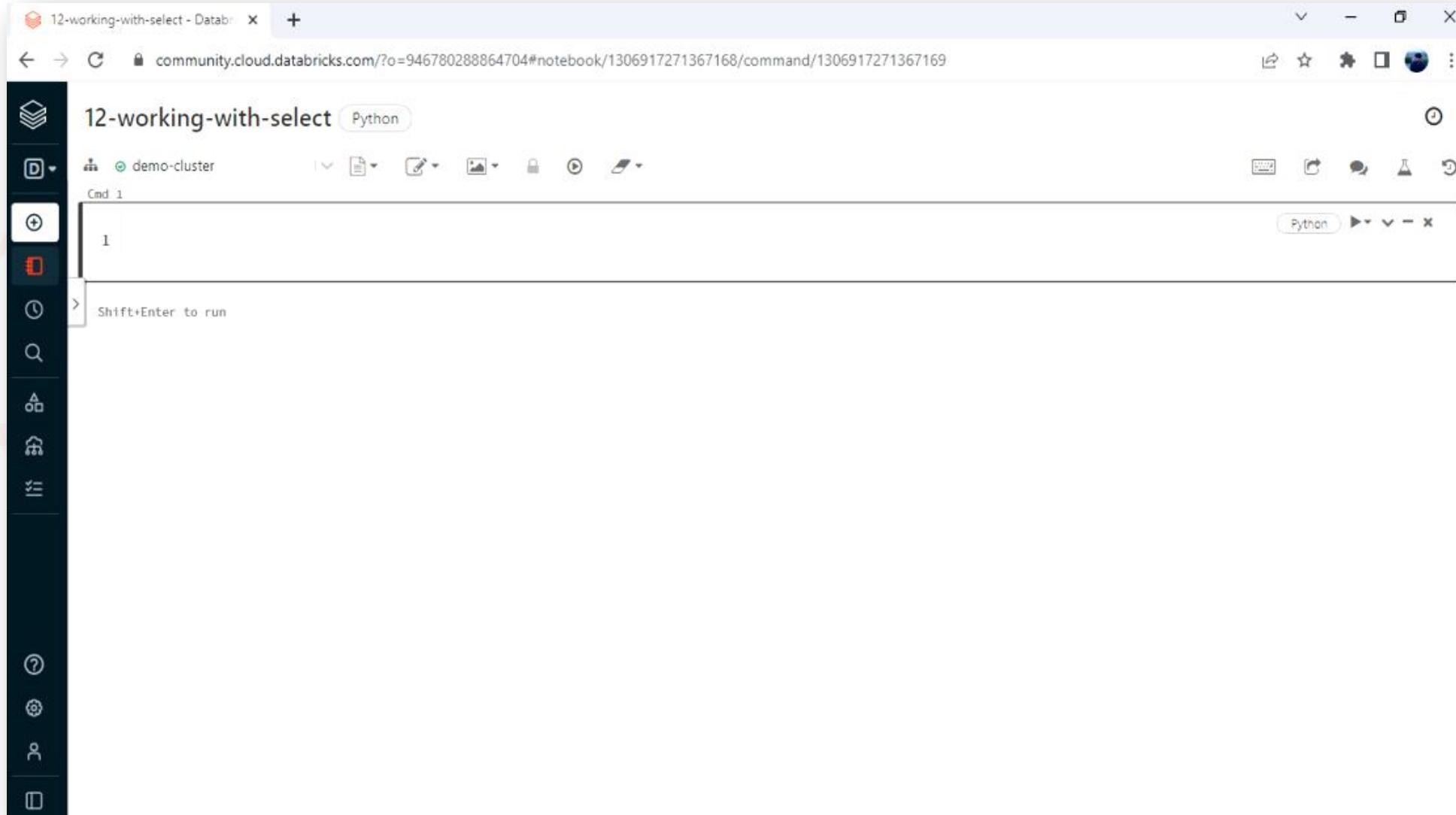


Working with Select and SelectExpr

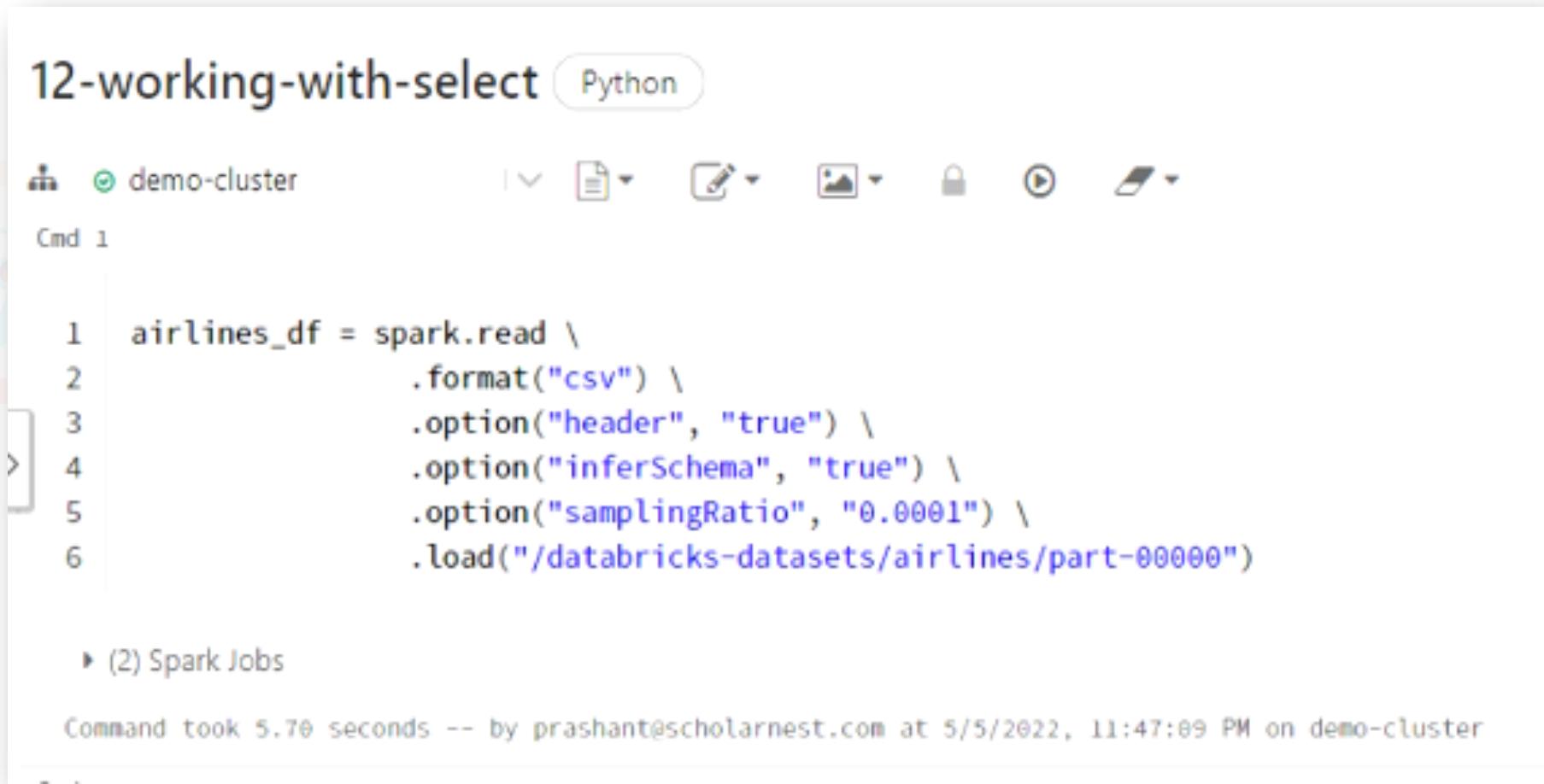
Spark gives you two methods to select one or more columns from your dataframe.

1. `select()`
2. `selectExpr()`

Go to your Databricks workspace and create a new notebook (**Reference - 12-working-with-select**)



We have created a Dataframe for using the select() method. I am reading one airline sample data from Databricks samples dataset and creating a Dataframe.



The screenshot shows a Databricks notebook interface with the title "12-working-with-select" and the language "Python". The notebook is connected to a "demo-cluster". The code in the first cell (Cmd 1) reads an airline dataset from a CSV file located at "/databricks-datasets/airlines/part-00000". The code is as follows:

```
1 airlines_df = spark.read \
2         .format("csv") \
3         .option("header", "true") \
4         .option("inferSchema", "true") \
5         .option("samplingRatio", "0.0001") \
6         .load("/databricks-datasets/airlines/part-00000")
```

Below the code, it shows "(2) Spark Jobs". At the bottom, the command took 5.70 seconds and was run by prashant@scholarnest.com on May 5, 2022, at 11:47:09 PM.

We have `airlines_df`, and we can start applying transformations. Let's assume this dataframe was a table. So you could run a `select *` query on this table like this - **SELECT * FROM table_name;**

The screenshot shows a Jupyter Notebook interface with a single code cell labeled "Cmd 2". The cell contains the command `display(airlines_df)`. Below the cell, the output is displayed as a truncated table of flight data. The table has 10 columns: Year, Month, DayofMonth, DayOfWeek, DepTime, CRSDepTime, ArrTime, CRSArrTime, UniqueCarrier, FlightNum, and TailNum. The first 7 rows of the table are shown, each corresponding to a different flight record from October 1987. The last row of the table indicates that only the first 1000 rows were displayed. The notebook header shows the title "12-working-with-select" and the language "Python". The status bar at the bottom of the screen shows the command took 0.93 seconds.

	Year	Month	DayofMonth	DayOfWeek	DepTime	CRSDepTime	ArrTime	CRSArrTime	UniqueCarrier	FlightNum	TailNum
1	1987	10	14	3	741	730	912	849	PS	1451	NA
2	1987	10	15	4	729	730	903	849	PS	1451	NA
3	1987	10	17	6	741	730	918	849	PS	1451	NA
4	1987	10	18	7	729	730	847	849	PS	1451	NA
5	1987	10	19	1	749	730	922	849	PS	1451	NA
6	1987	10	21	3	728	730	848	849	PS	1451	NA
7	1987	10	22	4	728	730	852	849	PS	1451	NA

Truncated results showing first 1000 rows:

Command took 0.93 seconds -- by prashant@scholarnest.com at 5/5/2022, 11:47:24 PM on demo-cluster

So I am selecting four columns here in this code. The select method takes a comma-separated list of column names and returns only those columns. You can give column names or column objects. The select method is designed to accept both, and you can also mix the argument types. You can give a column name or supply a column object in the same select expression.

```
Cmd 4

1 from pyspark.sql.functions import col, column
2
3 airlines_df.select("CRSDepTime", col("DepTime"), airlines_df.CRSArrTime, column("ArrTime")).show(5)

▶ (1) Spark Jobs
+-----+-----+-----+-----+
|CRSDepTime|DepTime|CRSArrTime|ArrTime|
+-----+-----+-----+-----+
|      730|    741|      849|    912|
|      730|    729|      849|    903|
|      730|    741|      849|    918|
|      730|    729|      849|    847|
|      730|    749|      849|    922|
+-----+-----+-----+-----+
only showing top 5 rows

Command took 0.85 seconds -- by prashant@scholarnest.com at 5/5/2022, 11:51:17 PM on demo-cluster
```

You can also give expressions that result in a column object as shown in the code below. The first argument is a column name. The second argument is also a column name, but I am wrapping it inside the `expr()` function. The `expr()` function evaluates an expression and returns a column object. We do not have any expression here. `DepTime` is just a column name. So I should not have wrapped it inside the `expr()` function. However, the `expr()` function doesn't complain. It will simply evaluate the column name to a column object. The third argument is an expression. I am deducting the scheduled departure time from the departure time to find the delay and renaming the result as `delay`.

```
Cmd 5

1 from pyspark.sql.functions import expr
2
3 airlines_df.select("CRSDepTime", expr("DepTime"), expr("DepTime - CRSDepTime as delay")).show(5)

▶ (1) Spark Jobs
+-----+-----+----+
|CRSDepTime|DepTime|delay|
+-----+-----+----+
|      730|     741|   11|
|      730|     729|  -1|
|      730|     741|   11|
|      730|     729|  -1|
|      730|     749|   19|
+-----+-----+----+
only showing top 5 rows

Command took 0.71 seconds -- by prashant@scholarnest.com at 5/5/2022, 11:53:05 PM on demo-cluster
```

The select() method takes only two types of arguments.

1. Column Name
2. Column Object

You can see some negative delays. I don't want them. I mean, the flight went early. But we are calculating delay here. So I want to see only positive delays and set all the negative values to 0.

How can we do it?

```
> Cmd 5

1 from pyspark.sql.functions import expr
2
3 airlines_df.select("CRSDepTime", expr("DepTime"), expr("DepTime - CRSDepTime as delay")).show(5)

▶ (1) Spark Jobs
+-----+-----+----+
|CRSDepTime|DepTime|delay|
+-----+-----+----+
|      730|     741|   11|
|      730|     729|   -1|
|      730|     741|   11|
|      730|     729|   -1|
|      730|     749|   19|
+-----+-----+----+
only showing top 5 rows

Command took 0.71 seconds -- by prashant@scholarnest.com at 5/5/2022, 11:53:05 PM on demo-cluster
```

Here is the code for that. I have the same SQL expression, and I evaluate it using the expr() function.

The expr() function will evaluate the expression to a column object, and it goes to the select() method.

```
> Cmd 6

1 airlines_df.select("CRSDepTime", "DepTime", expr("case when (DepTime - CRSDepTime) > 0 then (DepTime - CRSDepTime) else 0 end as delay")).show(5)

▶ (1) Spark Jobs

+-----+-----+-----+
|CRSDepTime|DepTime|delay|
+-----+-----+-----+
|      730|     741|    11|
|      730|     729|     0|
|      730|     741|    11|
|      730|     729|     0|
|      730|     749|    19|
+-----+-----+-----+
only showing top 5 rows

Command took 0.97 seconds -- by prashant@scholarnest.com at 5/6/2022, 12:01:08 AM on demo-cluster
```

Sometime, using the `expr()` function inside the `select()` looks inconvenient.

I mean, using expressions in the `select()` method is a common thing.

We often use it. But we have a shortcut to avoid the `expr()` function explicitly, and that is using the `selectExpr()` method. The `selectExpr()` is the same as the `select`, but it takes an expression.

The `select` method takes a column name or column object. However, the `selectExpr()` only takes SQL-like expressions.

```
> Cmd 7

1 airlines_df.selectExpr("CRSDepTime", "DepTime", "case when (DepTime - CRSDepTime) > 0 then (DepTime - CRSDepTime) else 0 end as delay").show(5)

▶ (1) Spark Jobs
+-----+-----+-----+
|CRSDepTime|DepTime|delay|
+-----+-----+-----+
|      730|     741|    11|
|      730|     729|     0|
|      730|     741|    11|
|      730|     729|     0|
|      730|     749|    19|
+-----+-----+-----+
only showing top 5 rows

Command took 0.49 seconds -- by prashant@scholarnest.com at 5/6/2022, 12:02:16 AM on demo-cluster
```

Here is another example.

Here third argument to the select() method is using a dataframe expression.

And the select() method takes it, because the expression ultimately evaluates to a column object. But if you try using selectExpr(), you will see an error.

```
> Cmd: 8  
1 from pyspark.sql.functions import col  
2  
3 some_v = 15  
4 airlines_df.select("CRSDepTime", "DepTime", col("DepTime") + some_v).show(5)
```

▶ (1) Spark Jobs

CRSDepTime	DepTime	(DepTime + 15)
730	741	756
730	729	744
730	741	756
730	729	744
730	749	764

only showing top 5 rows

Command took 0.70 seconds -- by prashant@scholarnest.com at 5/6/2022, 12:04:58 AM on demo-cluster

Here we are trying to use the `selectExpr()`, and you can see an error. So understand the difference between `select()` and `selectExpr()`. The `select` method takes a column name and a column object. You can use expressions as long as it evaluates a column object.

The `selectExpr()` only takes SQL-like expressions. You can also use column names because they evaluate as an expression.

The screenshot shows a Jupyter Notebook cell with the title "12-working-with-select" and the language "Python". The cell contains the following code:

```
1 from pyspark.sql.functions import col
2
3 some_v = 15
4 airlines_df.selectExpr("CRSDepTime", "DepTime", col("DepTime") + some_v).show(5)
```

The output of the code is:

	CRSDepTime	DepTime	DepTime + some_v
only showing top 5 rows			

Command took 0.54 seconds -- by prashant@scholarnest.com at 5/12/2022, 2:06:22 PM on demo-cluster

Below the code, an error message is displayed:

TypeError: Column is not iterable

Command took 0.43 seconds -- by prashant@scholarnest.com at 5/12/2022, 2:07:28 PM on demo-cluster



Thank You
ScholarNest Technologies Pvt Ltd.
www.scholarnest.com



ScholarNest

Spark Azure Databricks

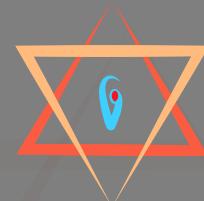
Databricks Spark Certification and beyond

Instructor: Prashant Kumar Pandey



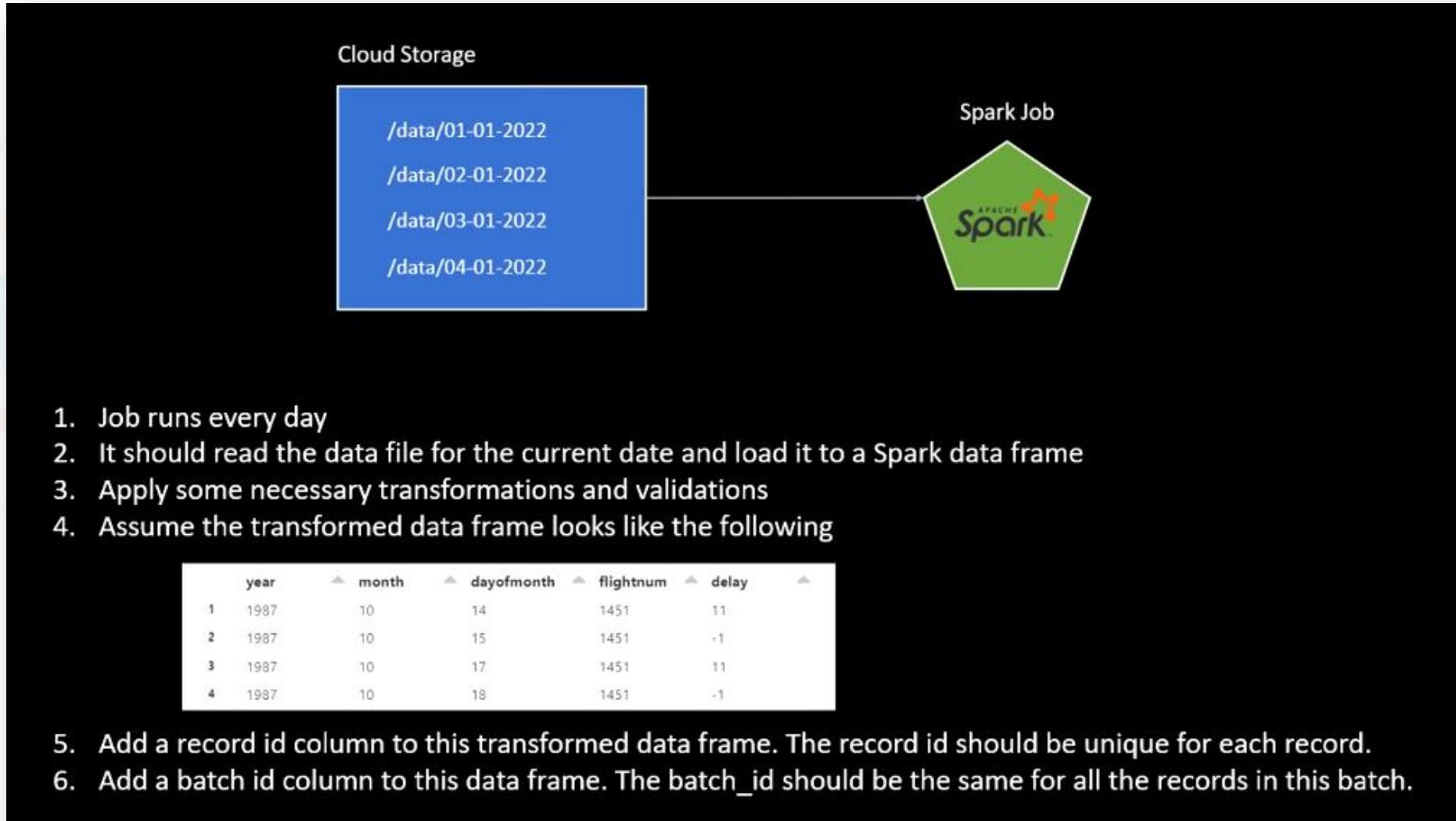
Absolute Beginner to Specialization in Apache Spark and Azure Databricks





Working with Literal Values

We have a simple requirement. You receive data in a cloud storage directory. The data comes every day in a date-wise directory. You are asked to create a Spark batch job to do the following.



The first requirement is to run the job every day.

And running a Job every day requires scheduling your Spark notebook or your Application file to run daily at a scheduled time.

We are not going into the details of scheduling for now.

I will cover Job scheduling in the latter part of the course.

So leave that requirement for now.

The second requirement is to read the data file for the current date.

There are many ways to achieve this requirement.

One easy method is to have date-wise directories.

Parameterize the load() method to read data from the current-date directory.

I will also leave this requirement and show you an example in the coming lecture.

The third requirement is to apply some transformations.

Let's leave this one also for now. We are learning transformations throughout the course.

And the focus of this lecture is on the fifth and the sixth requirement.

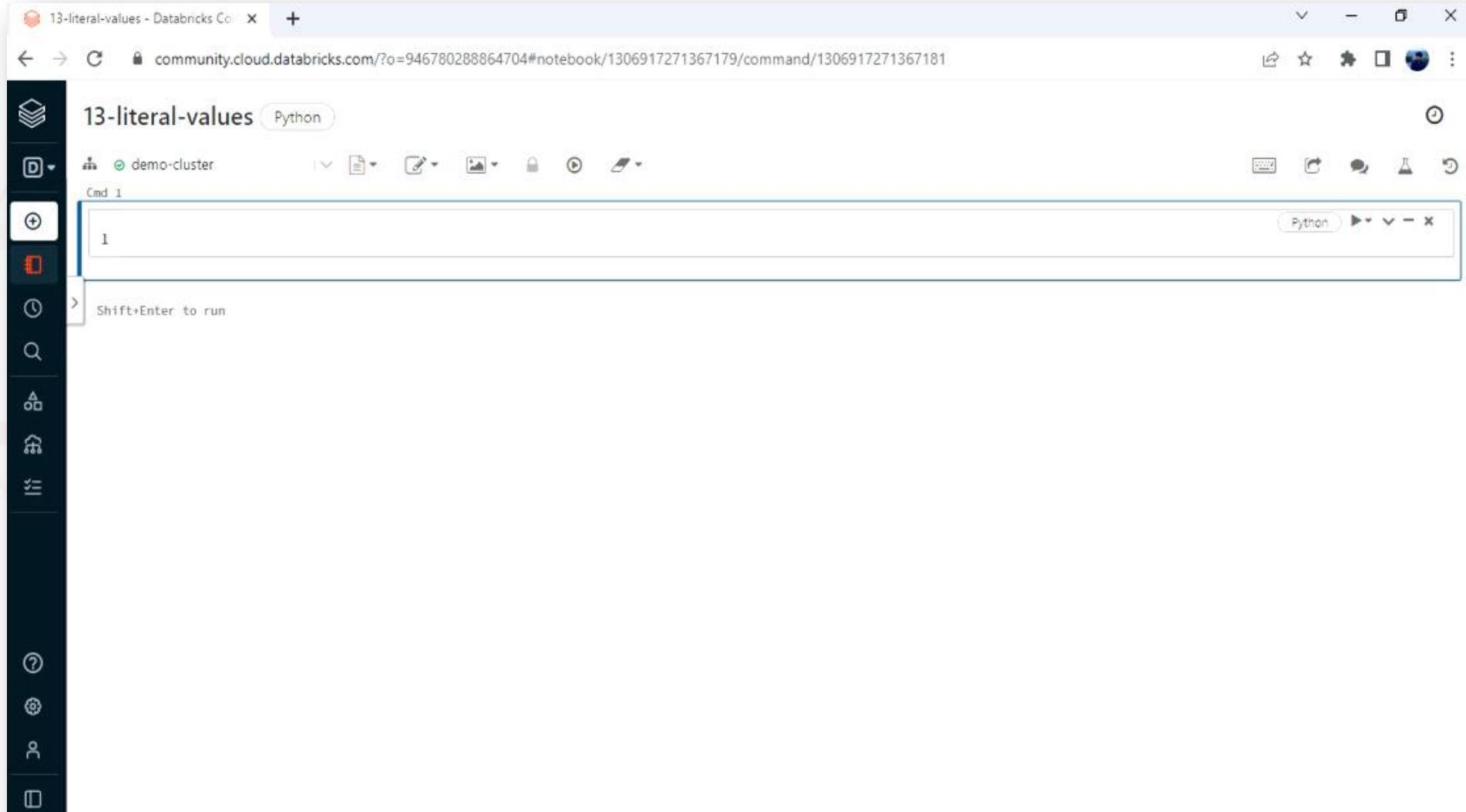
So let's leave this requirement also.

The fourth one shows you a sample dataframe.
So we will start from this requirement.
I will create this sample dataframe for you.
And you can learn how to implement the following requirements. Make sense?

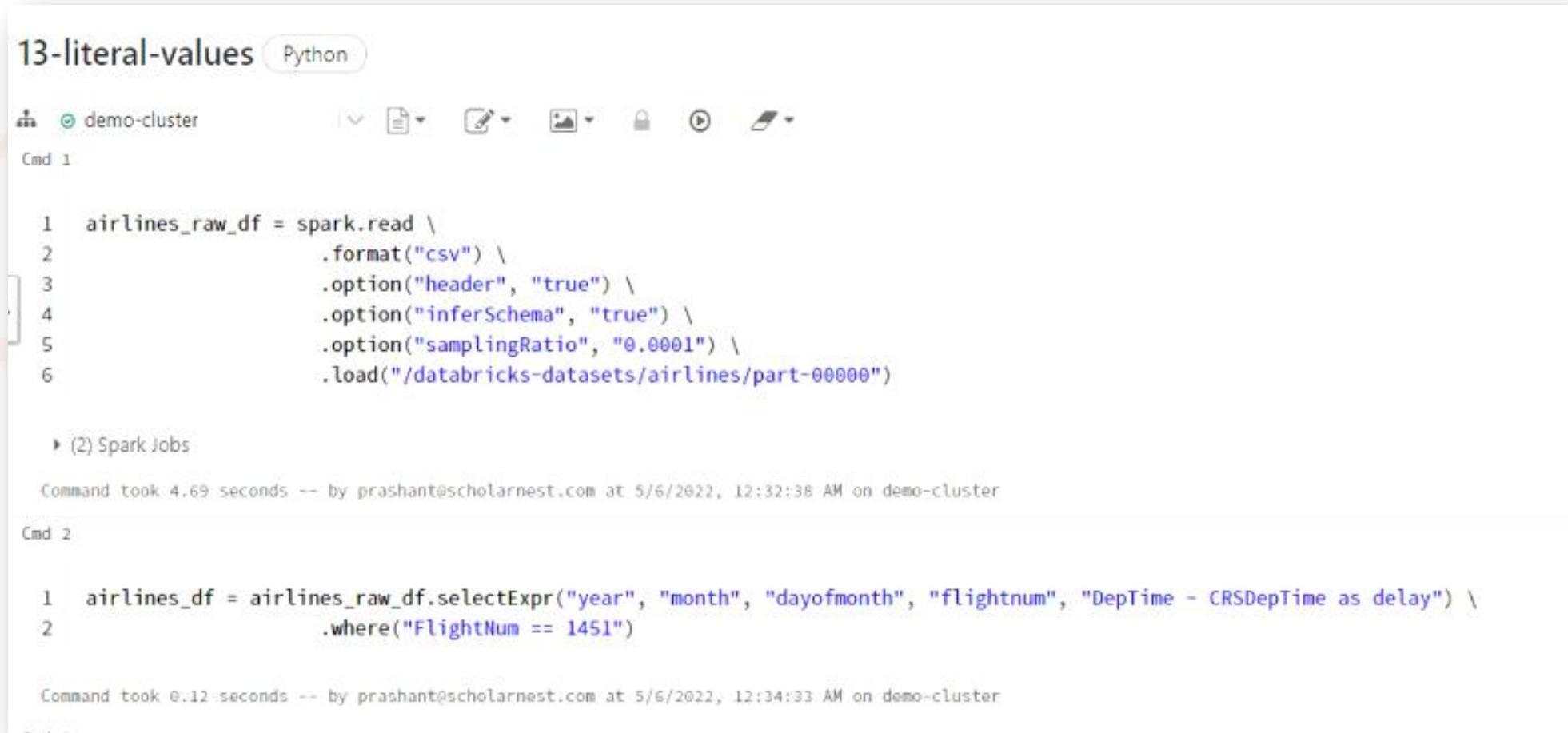
Great! So our primary focus is on the fifth and sixth items.
The fifth item wants us to add a unique record id for each record in the given dataframe.

Let's achieve this requirement, and then we will look into the last one.

Go to your Databricks workspace and create a new notebook. (**Reference : 13-literal-values**)



The first one is the same Dataframe that we created in our earlier lecture. We need to apply some transformations and prepare a Dataframe as per the example. So we transformed it and created a new Dataframe.



The screenshot shows a Databricks notebook interface with the title "13-literal-values" and a Python tab selected. The notebook contains two commands:

```
1 airlines_raw_df = spark.read \
2     .format("csv") \
3     .option("header", "true") \
4     .option("inferSchema", "true") \
5     .option("samplingRatio", "0.0001") \
6     .load("/databricks-datasets/airlines/part-00000")
```

After running Command 1, a message indicates "(2) Spark Jobs" were triggered, and the command took 4.69 seconds to execute.

```
▶ (2) Spark Jobs
Command took 4.69 seconds -- by prashant@scholarnest.com at 5/6/2022, 12:32:38 AM on demo-cluster
```

```
1 Cmd 2
2 airlines_df = airlines_raw_df.selectExpr("year", "month", "dayofmonth", "flightnum", "DepTime - CRSDepTime as delay") \
3     .where("FlightNum == 1451")
```

After running Command 2, a message indicates the command took 0.12 seconds to execute.

```
Command took 0.12 seconds -- by prashant@scholarnest.com at 5/6/2022, 12:34:33 AM on demo-cluster
```

We got the desired Dataframe. Now we have the following requirement:
Add a record id column to this transformed dataframe. The record id should be unique for each record.

```
Cmd 3

1 display(airlines_df)

▶ (3) Spark Jobs

+-----+-----+-----+-----+-----+
| year | month | dayofmonth | flightnum | delay |
+-----+-----+-----+-----+-----+
| 1    | 1987  | 10        | 14        | 1451    |
| 2    | 1987  | 10        | 15        | 1451    |
| 3    | 1987  | 10        | 17        | 1451    |
| 4    | 1987  | 10        | 18        | 1451    |
| 5    | 1987  | 10        | 19        | 1451    |
| 6    | 1987  | 10        | 21        | 1451    |
| 7    | 1987  | 10        | 22        | 1451    |
+-----+-----+-----+-----+-----+
Showing all 235 rows.

Command took 6.76 seconds -- by prashant@scholarnest.com at 5/6/2022, 12:34:50 AM on demo-cluster
```

How to do it? Can we use some function that generates a unique id for us?

We have two places to look for the available functions.

1. Built-in Spark SQL functions
2. Dataframe functions.

I checked both the places and identified two potential functions that we can use.

1. `uuid()`
2. `monotonically_increasing_id()`

The `uuid()` function is a built-in spark SQL function.

It generates a 32-character unique id.

We can use it to generate a UUID for each record.

The `monotonically_increasing_id()` is a dataframe function.

It generates a long number starting from 0.

We can use this function also to add a unique id for our records.

The monotonically_increasing_id() is a dataframe function. So you must import it. Once imported, you can use the withColumn() transformation to add a new column to your dataframe. The withColumn() takes two arguments. The first argument is the column name, and the second argument is the logic. The logic is as simple as using a function and generating a monotonically_increasing_id. You can see in the output there is the id column at the end, it starts with zero and goes on increasing by one.

13-literal-values Python

demo-cluster Command took 6.76 seconds -- by prashant@scholarnest.com at 5/6/2022, 12:34:50 AM on demo-cluster

Cmd 4

```
1 from pyspark.sql.functions import monotonically_increasing_id
2
3 airlines_df.withColumn("id", monotonically_increasing_id()) \
4     .show(5, False)
```

▶ (1) Spark Jobs

year	month	dayofmonth	flightnum	delay	id
1987	10	14	1451	11	0
1987	10	15	1451	-1	1
1987	10	17	1451	11	2
1987	10	18	1451	-1	3
1987	10	19	1451	19	4

only showing top 5 rows

Command took 0.54 seconds -- by prashant@scholarnest.com at 5/6/2022, 12:40:32 AM on demo-cluster

We have another alternative of using the UUID() built-in function. So I added another column and used the UUID() built-in function. Since the UUID() is a built-in function, I must use the expr() function to evaluate it. And we can see a unique id for each record. You can use any one of these depending upon your requirement.

13-literal-values Python

demo-cluster Command took 6.76 seconds -- by prashant@scholarnest.com at 5/6/2022, 12:34:50 AM on demo-cluster

Cmd 4

```
1 from pyspark.sql.functions import monotonically_increasing_id, expr
2
3 airlines_df.withColumn("id", monotonically_increasing_id()) \
4     .withColumn("RecordID", expr("uuid()")) \
5     .show(5, False)
```

▶ (1) Spark Jobs

year	month	dayofmonth	flightnum	delay	id	RecordID
1987	10	14	1451	11	0	bbfb4a69-2c6e-4f6e-a695-27edf83ccc74
1987	10	15	1451	-1	1	506c57dc-c092-4c5e-b555-5acb7905919f
1987	10	17	1451	11	2	b84dcefd-777c-40e2-94fa-c29b3e7bd769
1987	10	18	1451	-1	3	12fed1e8-7b7e-42a7-a699-5cde16d11e8e
1987	10	19	1451	19	4	5cf3a766-ac8e-40f7-9479-8dbba2e7f03f

only showing top 5 rows

Command took 0.52 seconds -- by prashant@scholarnest.com at 5/6/2022, 12:48:35 AM on demo-cluster

We want to add a batch id column to this dataframe. The batch_id should be the same for all the records in this batch. So it is like a hardcoded id for all the records.

How will you do it?

Can you even think of a SQL way of doing it?

Assume you have a table. You want to select all columns and add one extra column with a hardcoded value.

Here is a SQL expression to do it:

```
SELECT *, "batch-123456789" as batchID from table_name;
```

I wanted to add a new column to my Dataframe, so I used them with column transformation. The column name is BatchID, and the logic to create a column is a literal value. But if we try running it, we will get an error. Because the Dataframe doesn't understand the literal value. The second argument to the withColumn() must be an expression that evaluates a column object. This string is not an expression. I can wrap it with an expr() function, but even the expr() function will not evaluate it. Because the string here is a literal value. Expressions are created using column names, functions, and operators. And this string doesn't represent any column name. It is neither a function. It is a literal value, and Dataframe API cannot understand literal values until we explicitly tell it.

13-literal-values Python

demo-cluster Command took 6.76 seconds -- by prashant@scholarnest.com at 5/6/2022, 12:34:50 AM on demo-cluster

Cmd 4

```
1 from pyspark.sql.functions import monotonically_increasing_id, expr
2
3 airlines_df.withColumn("id", monotonically_increasing_id()) \
4     .withColumn("RecordID", expr("uuid()")) \
5     .withColumn("BatchID", "batch-123456789") \
6     .show(5, False)
```

TypeError: col should be Column

Command took 0.18 seconds -- by prashant@scholarnest.com at 5/6/2022, 12:49:01 AM on demo-cluster

I can import a Python UUID package and create a unique id to use the same example. Then I can use the python uid variable inside the lit() function.

So instead of hardcoding a manually created id, I can use a variable to automatically generate a unique id and use it as my batch id.

13-literal-values Python

demo-cluster Command took 6.76 seconds -- by prashant@scholarnest.com at 5/6/2022, 12:34:50 AM on demo-cluster

Cmd 4

```
1 from pyspark.sql.functions import monotonically_increasing_id, expr, lit
2
3 import uuid
4 uid = str(uuid.uuid4())
5
6 airlines_df.withColumn("id", monotonically_increasing_id()) \
    .withColumn("RecordID", expr("uuid()")) \
    .withColumn("BatchID", lit(uid)) \
    .show(5, False)
```

▶ (1) Spark Jobs

year	month	dayofmonth	flightnum	delay	id	RecordID	BatchID
1987	10	14	1451	11	0	5836314d-3702-4aa7-bd69-06ed0265ba97 932aed37-4837-4eaf-883c-540cc06c1132	
1987	10	15	1451	-1	1	aec77f55-2fcf-4c7b-bc65-bcf2a45f4bfd 932aed37-4837-4eaf-883c-540cc06c1132	
1987	10	17	1451	11	2	6dee1ebf-a0dd-4f43-b3d9-b2ed6b0260c9 932aed37-4837-4eaf-883c-540cc06c1132	
1987	10	18	1451	-1	3	0fc12fc1-d6e5-4280-9923-b88ee368e008 932aed37-4837-4eaf-883c-540cc06c1132	
1987	10	19	1451	19	4	8591b94f-ad50-4edb-ba25-91709408b2cc 932aed37-4837-4eaf-883c-540cc06c1132	

only showing top 5 rows

Command took 0.67 seconds -- by prashant@scholarnest.com at 5/6/2022, 12:51:01 AM on demo-cluster



Thank You
ScholarNest Technologies Pvt Ltd.
www.scholarnest.com



ScholarNest

Spark Azure Databricks

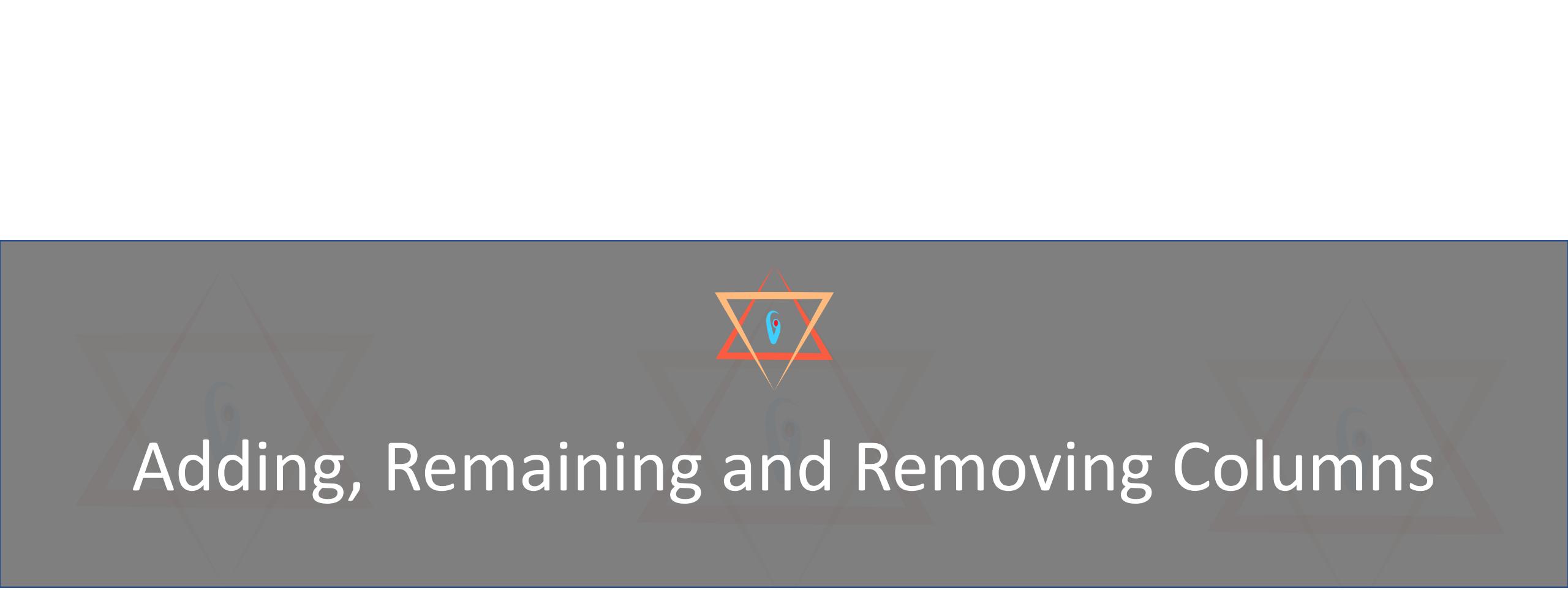
Databricks Spark Certification and beyond

Instructor: Prashant Kumar Pandey



Absolute Beginner to Specialization in Apache Spark and Azure Databricks



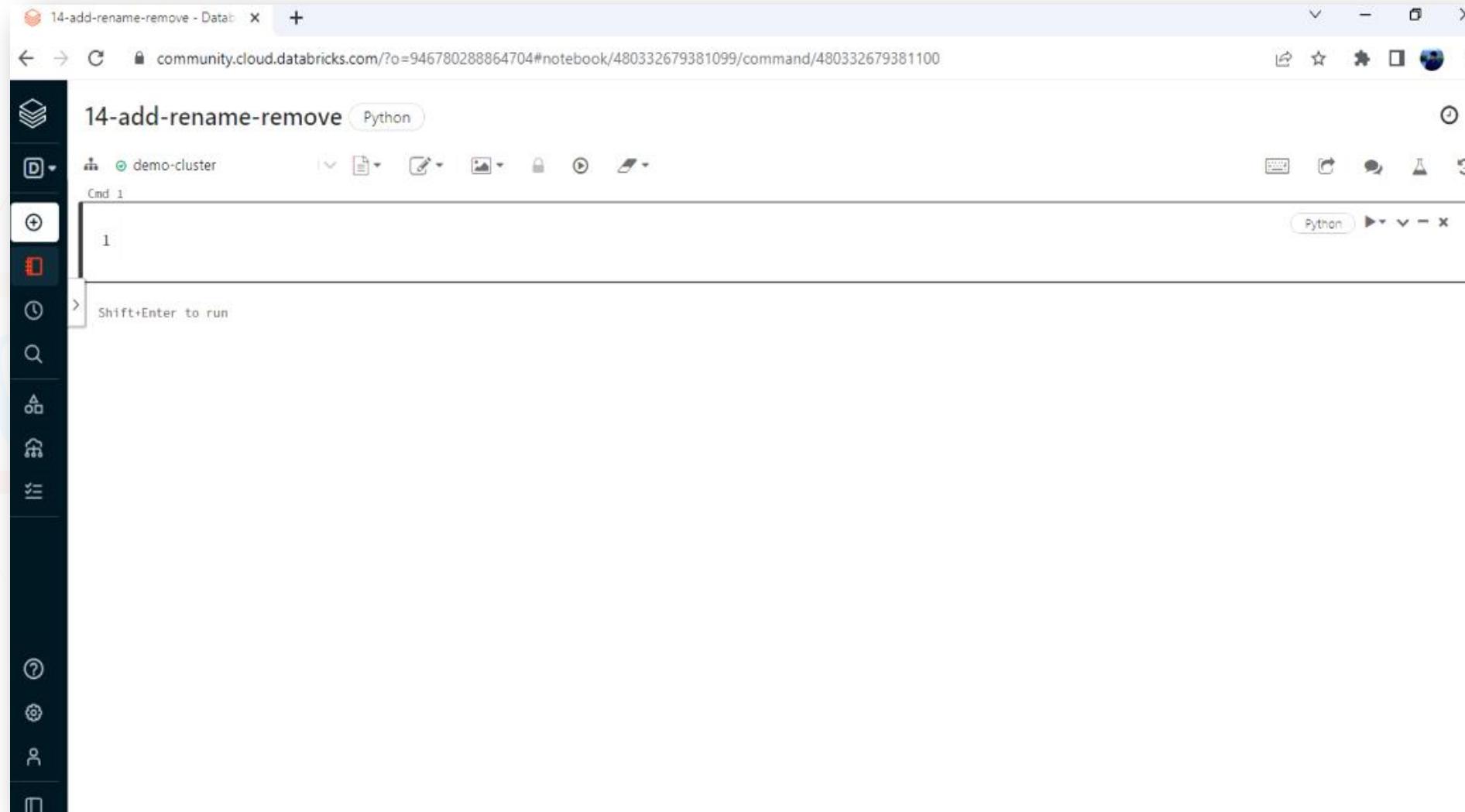


Adding, Remaining and Removing Columns

I want to talk about the following quickly in this chapter. The first thing that we want to cover is adding new columns to your dataframe. Spark dataframe gives you withColumn() transformation to do that. You can use withColumn() transformation to add a new column or transform an existing column. We have already used it in many places, and I hope you have already learned it. But let me give you a quick example once again.

1. Adding new columns to your data frame
 - withColumn()
2. Renaming Existing Columns in your data frame
 - withColumnRenamed()
3. Dropping columns from your data frame
 - drop()

Go to your Databricks workspace and create a new notebook (**Reference : 14-add-rename-remove**).



I am creating a data list and using it to define a data frame. You already learned this in an earlier lecture. Now I want to use a new column for this data frame. So I can use the `withColumn()` transformation.

14-add-rename-remove Python

demo-cluster

Cmd 1

```
1 data_list = [(100, "Prashant", 45, 45000),
 2                 (101, "Tarun", 36, 33000),
 3                 (102, "David", 48, 28000)]
4
5 test_df = spark.createDataFrame(data_list).toDF("id", "name", "age", "salary")
6 display(test_df)
```

(3) Spark Jobs

	id	name	age	salary
1	100	Prashant	45	45000
2	101	Tarun	36	33000
3	102	David	48	28000

Showing all 3 rows.

Command took 6.04 seconds -- by prashant@scholarnest.com at 5/9/2022, 2:45:04 PM on demo-cluster

So you can use the `withColumn()` method to add a new field to your dataframe. The `withColumn()` method takes two arguments. The first argument is the column name and the second column is an expression to determine the value of the new columns. The `withColumn()` can add new columns or transform existing columns. The increment column doesn't exist in this dataframe, so the `withColumn()` will add it. However, if the column already exists, the `withColumn()` will simply transform it.

```
> Cmd 2

1 from pyspark.sql.functions import expr
2
3 test_df1 = test_df.withColumn("increment", expr("salary * 10/100"))
4 display(test_df1)

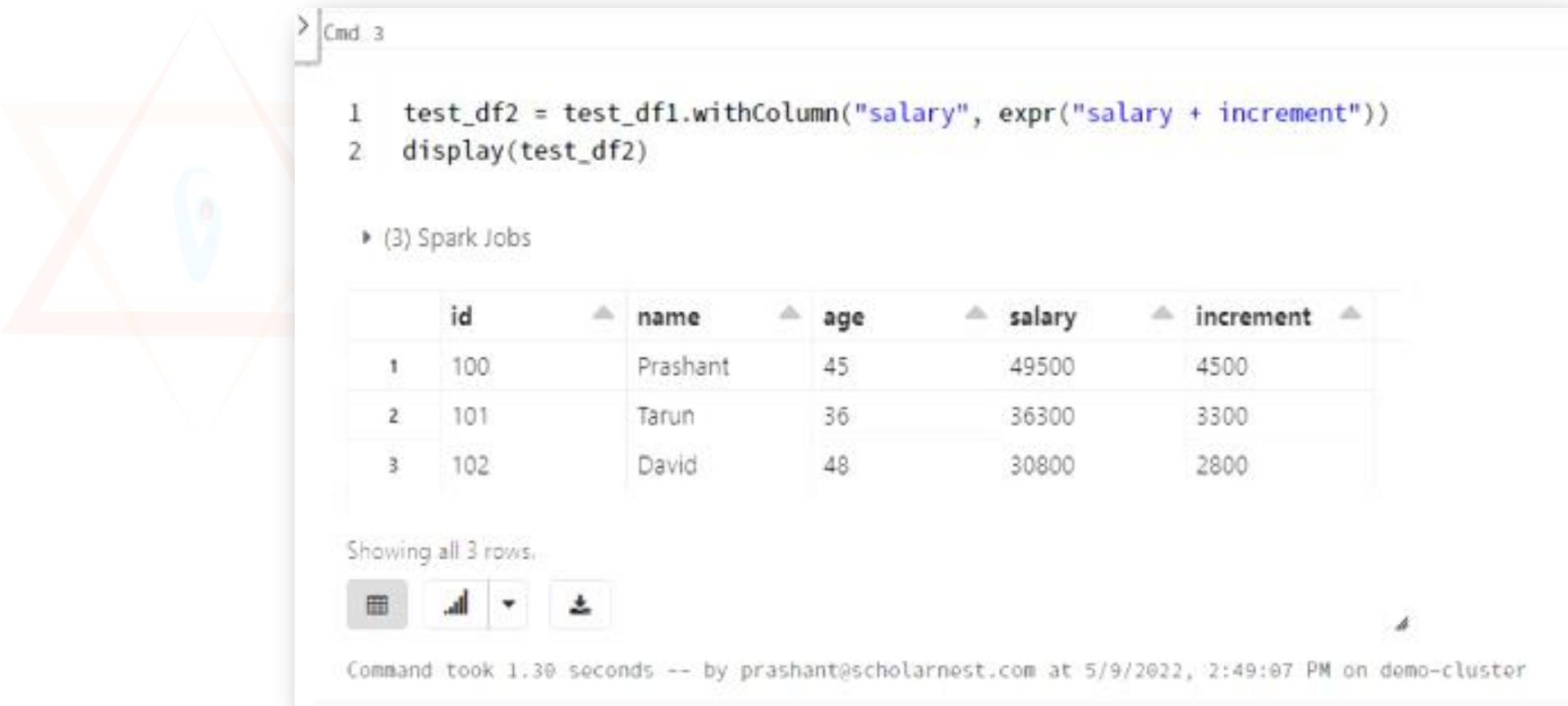
▶ (3) Spark Jobs

+---+---+---+---+---+
| id | name | age | salary | increment |
+---+---+---+---+---+
| 1  | Prashant | 45 | 45000 | 4500 |
| 2  | Tarun    | 36 | 33000 | 3300 |
| 3  | David    | 48 | 28000 | 2800 |
+---+---+---+---+---+

Showing all 3 rows.
[CSV, JSON, PDF, Copy]
```

Command took 1.57 seconds -- by prashant@scholarnest.com at 5/9/2022, 2:47:13 PM on demo-cluster

The salary column already existed, so the `withColumn()` updates the same column. The `withColumn()` is super convenient. However, it has got some problems. Each `withColumn()` method introduces a new projection in the SQL execution plan. Therefore, calling it multiple times can generate big SQL execution plans, which can cause performance issues and even `StackOverflowException`. To avoid this, we prefer to use the `select()` expression and avoid using the long list of `withColumn()` method.



```
> Cmd 3

1 test_df2 = test_df1.withColumn("salary", expr("salary + increment"))
2 display(test_df2)

▶ (3) Spark Jobs



|   | <b>id</b> | <b>name</b> | <b>age</b> | <b>salary</b> | <b>increment</b> |
|---|-----------|-------------|------------|---------------|------------------|
| 1 | 100       | Prashant    | 45         | 49500         | 4500             |
| 2 | 101       | Tarun       | 36         | 36300         | 3300             |
| 3 | 102       | David       | 48         | 30800         | 2800             |



Showing all 3 rows.

Command took 1.30 seconds -- by prashant@scholarnest.com at 5/9/2022, 2:49:07 PM on demo-cluster
```

Now look at both codes and compare them.

I am creating `data_list1` and then creating a dataframe using the same.

Then I am using the `withColumn()` method to add an increment column.

Then I am again using a `withColumn()` to change the salary column value.

Finally, I will show `my_df1`.

14-add-rename-remove Python

demo-cluster Cmd 4

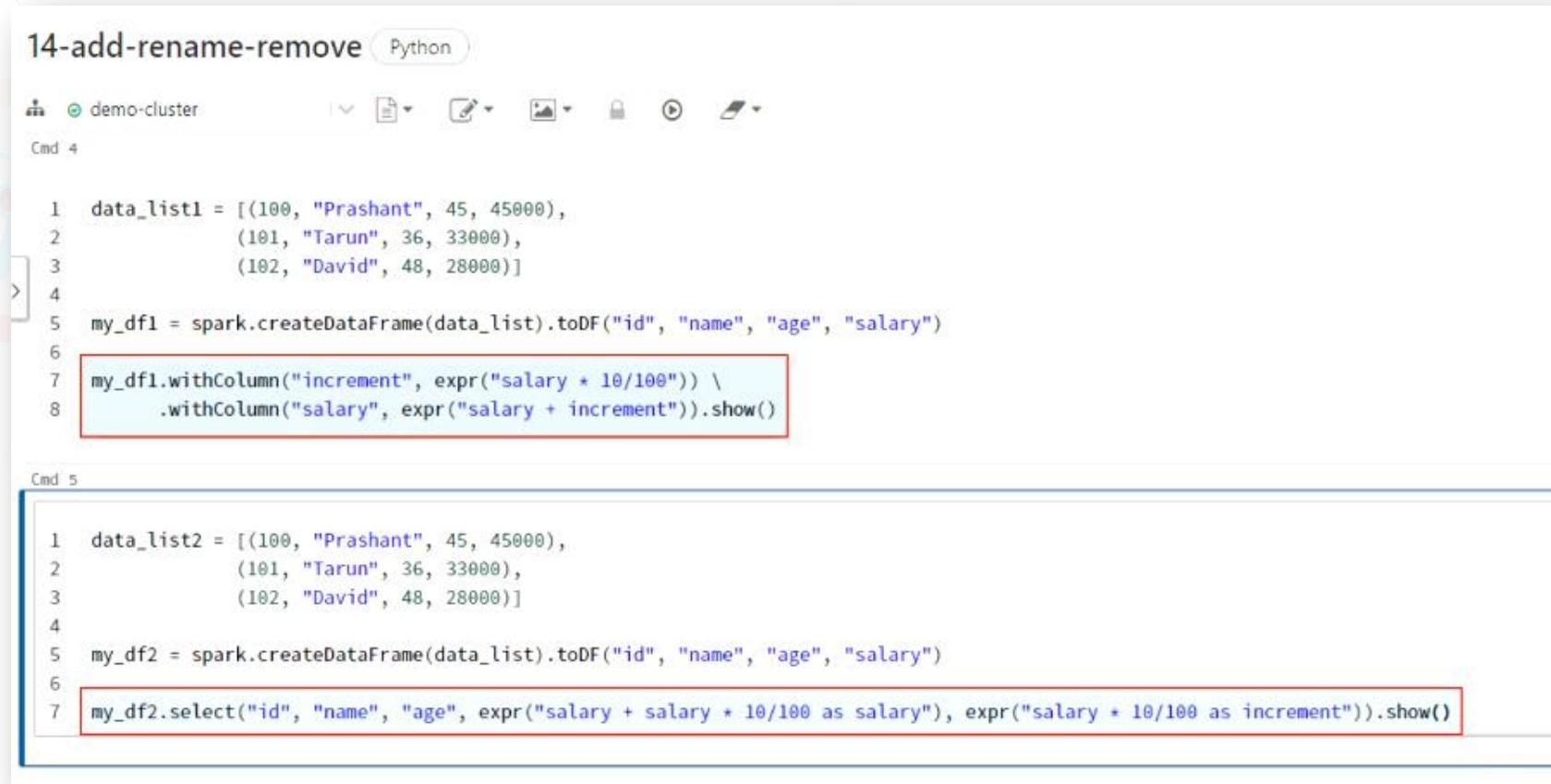
```
1 data_list1 = [(100, "Prashant", 45, 45000),
2                 (101, "Tarun", 36, 33000),
3                 (102, "David", 48, 28000)]
4
5 my_df1 = spark.createDataFrame(data_list).toDF("id", "name", "age", "salary")
6
7 my_df1.withColumn("increment", expr("salary * 10/100")) \
8     .withColumn("salary", expr("salary + increment")).show()
```

Cmd 5

```
1 data_list2 = [(100, "Prashant", 45, 45000),
2                 (101, "Tarun", 36, 33000),
3                 (102, "David", 48, 28000)]
4
5 my_df2 = spark.createDataFrame(data_list).toDF("id", "name", "age", "salary")
6
7 my_df2.select("id", "name", "age", expr("salary + salary * 10/100 as salary"), expr("salary * 10/100 as increment")).show()
```

Now, look at the next cell.

I am doing the same thing here but taking a different approach. This time, I am creating data_list2, which is precisely the same as data_list1. Then I am using data_list2 to create my_df2. The my_df2 is also precisely the same as my_df1. Finally, I am using a select() method. The select method selects id, name, and age without any change. But I am using an expression to change the salary column value. Then I use another expression to add an increment column to my dataframe.



The screenshot shows a Jupyter Notebook interface with two code cells. The first cell, titled '14-add-rename-remove' and labeled 'Python', contains the following code:

```
1 data_list1 = [(100, "Prashant", 45, 45000),
2                 (101, "Tarun", 36, 33000),
3                 (102, "David", 48, 28000)]
4
5 my_df1 = spark.createDataFrame(data_list).toDF("id", "name", "age", "salary")
6
7 my_df1.withColumn("increment", expr("salary * 10/100")) \
8     .withColumn("salary", expr("salary + increment")).show()
```

The last two lines of this cell are highlighted with a red box.

The second cell, labeled 'Cmd 5', contains the following code:

```
1 data_list2 = [(100, "Prashant", 45, 45000),
2                 (101, "Tarun", 36, 33000),
3                 (102, "David", 48, 28000)]
4
5 my_df2 = spark.createDataFrame(data_list).toDF("id", "name", "age", "salary")
6
7 my_df2.select("id", "name", "age", expr("salary + salary * 10/100 as salary"), expr("salary * 10/100 as increment")).show()
```

The last line of this cell is highlighted with a red box.

Now compare the two `withColumn()` expressions in the first cell with the `select()` method in the second cell. Both the codes are doing the same thing. However, they take a different approach. The first approach is using a chain of two `withColumn()` methods. The second approach is using a single `select()` method. Coding the first approach is neat and convenient. But coding the second approach is a bit clumsy.

14-add-rename-remove Python

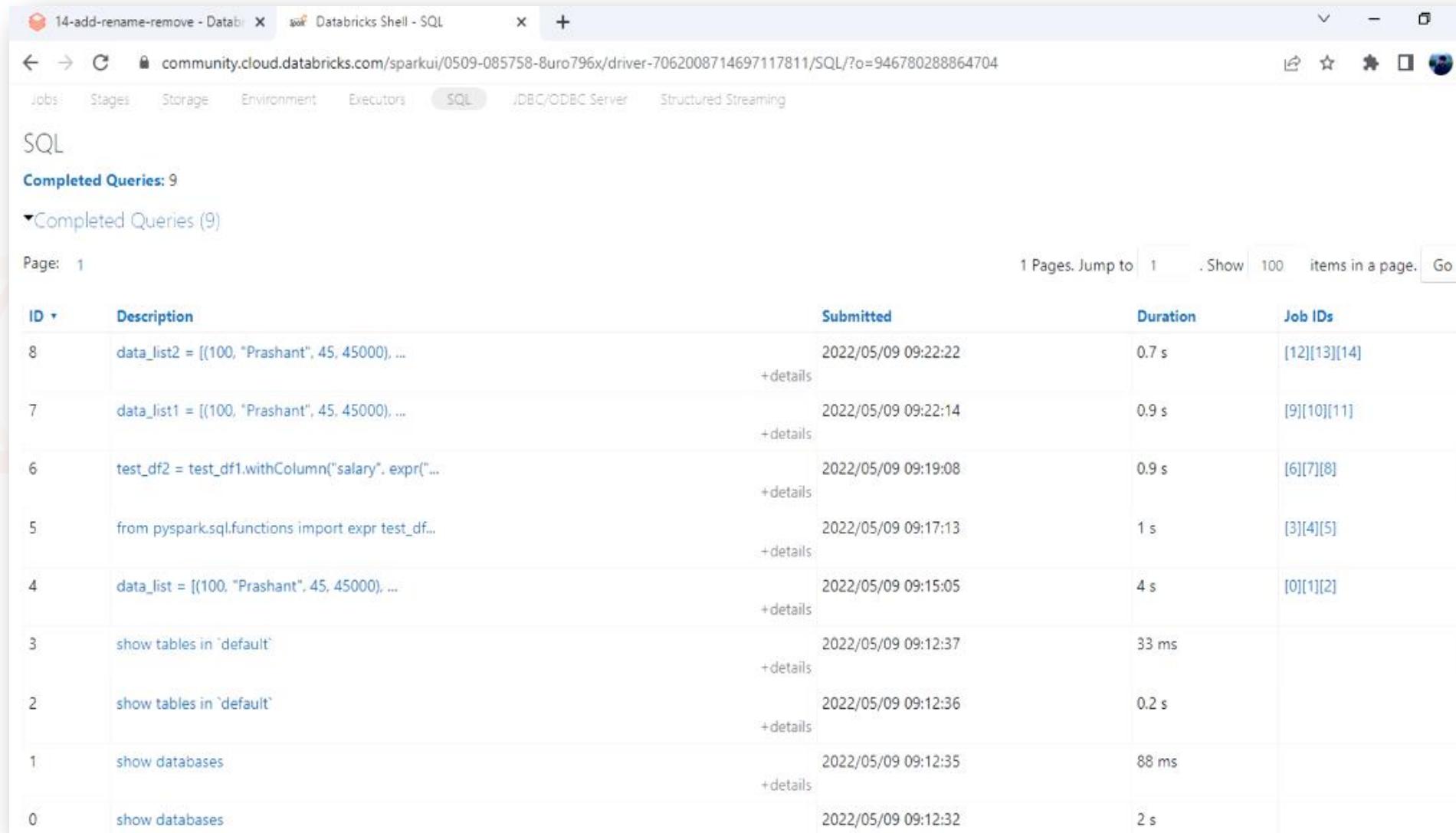
demo-cluster Cmd 4

```
1 data_list1 = [(100, "Prashant", 45, 45000),
2                 (101, "Tarun", 36, 33000),
3                 (102, "David", 48, 28000)]
4
5 my_df1 = spark.createDataFrame(data_list).toDF("id", "name", "age", "salary")
6
7 my_df1.withColumn("increment", expr("salary * 10/100")) \
8     .withColumn("salary", expr("salary + increment")).show()
```

Cmd 5

```
1 data_list2 = [(100, "Prashant", 45, 45000),
2                 (101, "Tarun", 36, 33000),
3                 (102, "David", 48, 28000)]
4
5 my_df2 = spark.createDataFrame(data_list).toDF("id", "name", "age", "salary")
6
7 my_df2.select("id", "name", "age", expr("salary + salary * 10/100 as salary"), expr("salary * 10/100 as increment")).show()
```

To check which one performs better. Run both the cells shown in the last slide, and check for the Spark UI, then go to the SQL tab. Open ID 8 and 7 in new tab.



The screenshot shows the Databricks SQL tab interface. 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, the URL is community.cloud.databricks.com/sparkui/0509-085758-8uro796x/driver-7062008714697117811/SQL/?o=946780288864704. The main area is titled "SQL" and displays "Completed Queries: 9". A dropdown menu "Completed Queries (9)" is open. Below it, there is a table with the following data:

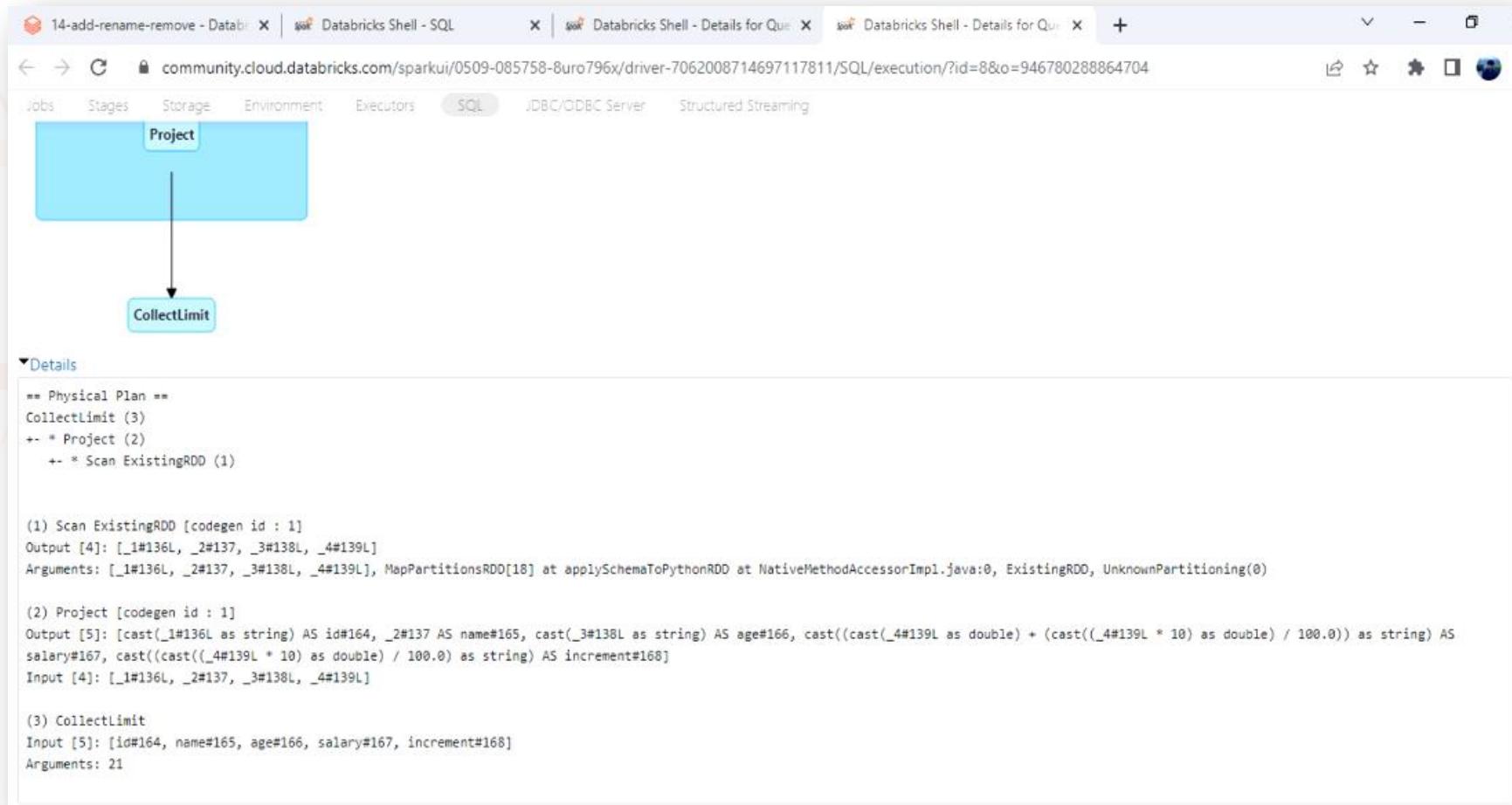
ID	Description	Submitted	Duration	Job IDs
8	data_list2 = [(100, "Prashant", 45, 45000), ...]	2022/05/09 09:22:22 +details	0.7 s	[12][13][14]
7	data_list1 = [(100, "Prashant", 45, 45000), ...]	2022/05/09 09:22:14 +details	0.9 s	[9][10][11]
6	test_df2 = test_df1.withColumn("salary", expr("..."))	2022/05/09 09:19:08 +details	0.9 s	[6][7][8]
5	from pyspark.sql.functions import expr test_df...	2022/05/09 09:17:13 +details	1 s	[3][4][5]
4	data_list = [(100, "Prashant", 45, 45000), ...]	2022/05/09 09:15:05 +details	4 s	[0][1][2]
3	show tables in `default`	2022/05/09 09:12:37 +details	33 ms	
2	show tables in `default`	2022/05/09 09:12:36 +details	0.2 s	
1	show databases	2022/05/09 09:12:35 +details	88 ms	
0	show databases	2022/05/09 09:12:32	2 s	

Come to the first one and check for the physical plan. The plan shows four steps, and we have two project operations in the SQL plan. So this one is the plan for the two consecutive withColumn() methods.

The screenshot shows the Databricks Shell interface with the 'SQL' tab selected. A blue arrow points down to the 'CollectLimit' step in the physical plan details. The physical plan consists of four steps:

- (1) Scan ExistingRDD [codegen id : 1]**
Output [4]: [_1#87L, _2#88, _3#89L, _4#90L]
Arguments: [_1#87L, _2#88, _3#89L, _4#90L], MapPartitionsRDD[12] at applySchemaToPythonRDD at NativeMethodAccessorImpl.java:0, ExistingRDD, UnknownPartitioning(0)
- (2) Project [codegen id : 1]**
Output [5]: [_1#87L AS id#95L, _2#88 AS name#96, _3#89L AS age#97L, _4#90L AS salary#98L, (cast(_4#90L * 10) as double) / 100.0 AS increment#103]
Input [4]: [_1#87L, _2#88, _3#89L, _4#90L]
- (3) Project [codegen id : 1]**
Output [5]: [cast(id#95L as string) AS id#120, name#96, cast(age#97L as string) AS age#122, cast((cast(salary#98L as double) + increment#103) as string) AS salary#123, cast(increment#103 as string) AS increment#124]
Input [5]: [id#95L, name#96, age#97L, salary#98L, increment#103]
- (4) CollectLimit**
Input [5]: [id#120, name#96, age#122, salary#123, increment#124]
Arguments: 21

Come to the second one and check for the physical plan. This plan shows only three steps and only one project operation. So the select() method worked faster. Every withColumn() adds one project operation in the SQL plan and slows down the SQL. So we recommend combining multiple withColumn() methods with a single select() method. This approach will cut down extra projections and speed up your SQL.



Spark dataframe gives you the `withColumnRenamed()` method to rename an existing column. Here is an example shown below. I am using a chain of two `withColumnRenamed()` methods to rename two columns. You should note two things here.

1. The `withColumnRenamed()` is a no-op if the schema doesn't contain the given column name.
2. The `withColumnRenamed()` will also add one projection to your SQL plan. So chaining multiple `withColumnRenamed()` is not recommended.

14-add-rename-remove Python

demo-cluster

Command took 1.06 seconds -- by prashant@scholarnest.com at 5/9/2022, 2:52:21 PM on demo-cluster

Cmd 6

```
1 test_df2.withColumnRenamed("increment", "salary_increment") \
2     .withColumnRenamed("salary", "incremented_salary").show()
```

▶ (3) Spark Jobs

id	name	age	incremented_salary	salary_increment
100	Prashant	45	49500.0	4500.0
101	Tarun	36	36300.0	3300.0
102	David	48	30800.0	2800.0

Command took 0.62 seconds -- by prashant@scholarnest.com at 5/9/2022, 2:59:02 PM on demo-cluster

Here are some facts about column naming.

1. Dataframe Column names are case insensitive.
 2. You might be able to give a wizard name with spaces and use keywords without any problem.
But you will face problems when using those names in expressions.

In the example shown below, I renamed the salary column to a wizard name.

This new name has upper and lower case letters. We also have a keyword. The as is a keyword. I also have an arithmetic symbol in the column name. Such names are complete nonsense. But you can do it.

14-add-rename-remove Python

demo-cluster | ▾

Command took 0.5 seconds -- by prashant@scholarnest.com at 5/9/2022, 2:59:48 PM on demo-cluster

Cmd 8

```
1 test_df3 = test_df2.withColumnRenamed("salary", "SaLArY as INcrement + sALArY")
2 display(test_df3)
```

▶ (3) Spark Jobs

	id	name	age	SaLArY as INcrement + sALArY	increment
1	100	Prashant	45	49500	4500
2	101	Tarun	36	36300	3300
3	102	David	48	30800	2800

Showing all 3 rows.

grid chart download

Command took 0.51 seconds -- by prashant@scholarnest.com at 5/9/2022, 3:01:36 PM on demo-cluster

So you won't have any problem creating crappy column names. But you face a problem when you want to use them in an expression. In the code shown below, I want to subtract increment from the "salary as increment + salary." I am using all small case letters in the column name. And that's fine because Spark Dataframe column names are not case-sensitive. But we get an error.

```
Cmd 9

1 test_df3.withColumn("Salary as increment + salary", expr("salary as increment + salary - increment")).show()

ParseException:
mismatched input '+' expecting {, '-'}(line 1, pos 20)

== SQL ==
salary as increment + salary - increment
-----^

Command took 0.20 seconds -- by prashant@scholarnest.com at 5/9/2022, 3:03:51 PM on demo-cluster
```

You should not name your columns in such away.

The recommendation is to follow a standard naming convention and avoid giving such weird names.

But if you do, you can still use them in the expression. Enclose your column names in a pair of backtick characters. You will find the backtick character below the escape key on your keyboard. And you can see that it works.

```
> Cmd 9

1 test_df3.withColumn("Salary as increment + salary", expr("`salary as increment + salary` - increment")).show()

▶ (3) Spark Jobs

+---+-----+---+-----+
| id|    name|age|Salary as increment + salary|increment|
+---+-----+---+-----+
|100|Prashant| 45|      45000.0|   4500.0|
|101| Tarun| 36|      33000.0|   3300.0|
|102| David| 48|      28000.0|   2800.0|
+---+-----+---+-----+

Command took 0.95 seconds -- by prashant@scholarnest.com at 5/9/2022, 3:04:52 PM on demo-cluster
```

So the point is straight.

Do not create column names with spaces, keywords, and symbols.

But you may still see such column names when you read data from external systems and sources.

Spark allows you to enclose such column names in a pair of backtick and use them in your expressions.

But you should rename such columns and give them a meaningful standard name.

To delete or remove a column from your Dataframe. Spark gives you a drop() method. Here in this example, I want to drop two columns from the test_df3 Dataframe. So I can use the drop() method and pass a comma-separated list of columns. It will drop all the given columns and give you a new Dataframe.

```
> Cmd 10

1 test_df3.drop("salary as increment + salary", "increment").show()

▶ (3) Spark Jobs

+---+---+
| id| name| age|
+---+---+
| 100| Prashant| 45|
| 101| Tarun| 36|
| 102| David| 48|
+---+---+


Command took 0.74 seconds -- by prashant@scholarnest.com at 5/9/2022, 3:06:05 PM on demo-cluster
```



Thank You
ScholarNest Technologies Pvt Ltd.
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



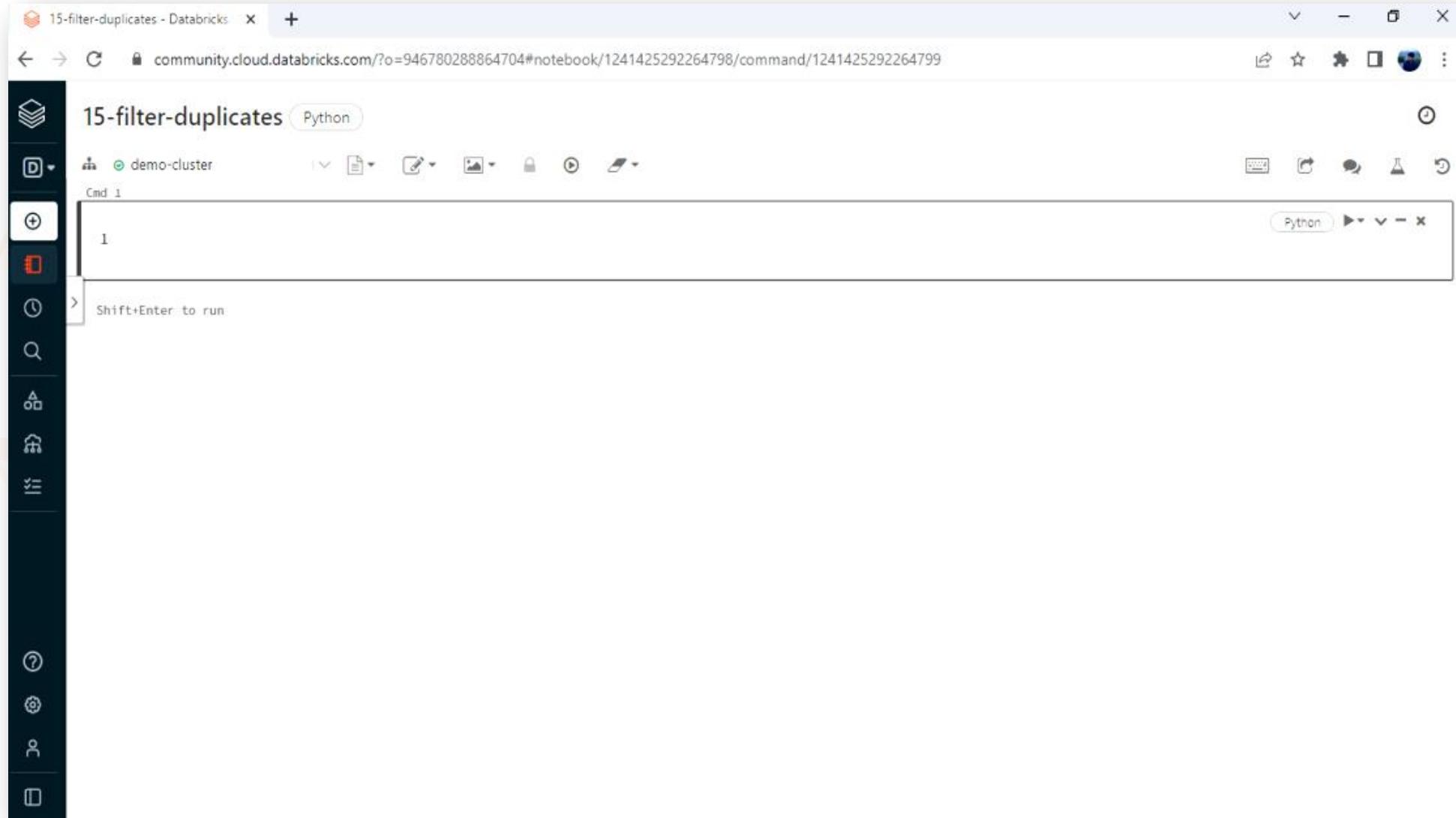


Filtering Records and Finding Unique Rows

In this chapter, we will talk about the following:

1. How to filter records from your Spark Dataframe
2. How to eliminate duplicates from your Dataframe

Go to your Databricks workspace and create a new notebook (**Reference : 15-filter-duplicates**)



Here we have a Dataframe. Now let's assume we want to filter records where the source column starts with Mum.

15-filter-duplicates Python

```
demo-cluster
8
9 df = spark.createDataFrame(data_list) \
10     .toDF("id", "source", "destination", "distance")
11
12 display(df)
```

(3) Spark Jobs

	id	source	destination	distance
1	101	Mumbai	Goa	587
2	102	Mumbai	Bangalore	985
3	102	Mumbai	Bangalore	985
4	103	Delhi	Chennai	2208
5	104	Delhi	Chennai	2208
6	105	Bangalore	Kolkata	1868
7	105	Bangalore	Kolkata	1865

Showing all 7 rows.

Command took 6.19 seconds -- by prashant@scholarnest.com at 5/11/2022, 8:33:04 PM on demo-cluster

To filter records where the source column starts with Mum, we can use the following SQL command.

```
SELECT * FROM table_name  
WHERE source like 'Mum%'
```

How to do it in a dataframe?

Spark gives you two methods.

1. `filter()`
2. `where()`

The `filter` method allows you to filter records in the dataframe.

And the `where()` method is an alias of the `filter()` method.

So both are the same.

The filter condition is the same as the condition in a SQL expression. Then, I am chaining one more condition using the where method. There is no difference between filter() and where() because the where() method is just an alias of the filter() method.

```
> Cmd 2

1 df.filter("source like 'Mum%'") \
2     .where("destination like 'Ban%'") \
3     .show()

▶ (3) Spark Jobs

+---+-----+-----+
| id|source|destination|distance|
+---+-----+-----+
| 102|Mumbai| Bangalore|    985|
| 102|Mumbai| Bangalore|    985|
+---+-----+-----+

Command took 0.79 seconds -- by prashant@scholarnest.com at 5/11/2022, 8:51:47 PM on demo-cluster
```

I used SQL conditions in the filter and where methods. And that's the most convenient method. But the filter method also allows you to specify conditions using functions and column API. This approach to writing expression is complex and confusing for SQL developers.

And there is no harm in taking the SQL approach and defining your expressions.

But learning the functional approach is good to have even though you rarely use it in practice. So this code takes the source column and applies column API to add the like() condition.

```
? Cmd: 3

1 from pyspark.sql.functions import col
2
3 df.filter(col("source").like("Mum%")) \
4 .where(col("destination").like("Ban%")) \
5 .show()

▶ (3) Spark Jobs
+---+-----+-----+
| id|source|destination|distance|
+---+-----+-----+
| 102|Mumbai| Bangalore|    985|
| 102|Mumbai| Bangalore|    985|
+---+-----+-----+

Command took 0.58 seconds -- by prashant@scholarnest.com at 5/9/2022, 3:52:37 PM on demo-cluster
```

We can also use the previous example like this. I am explicitly using the AND operator between two conditions.

And this approach is more intuitive.

So we recommend using conditions explicitly.

Cmd 4

```
1 df.filter("source like 'Mum%' and destination like 'Ban%'") \
2     .show()
```

▶ (3) Spark Jobs

id	source	destination	distance
102	Mumbai	Bangalore	985
102	Mumbai	Bangalore	985

Command took 0.59 seconds -- by prashant@scholarnest.com at 5/9/2022, 3:54:13 PM on demo-cluster

Column APIs also allow using the and condition. The syntax is a little different as shown below. So we use a single ampersand between two conditions.

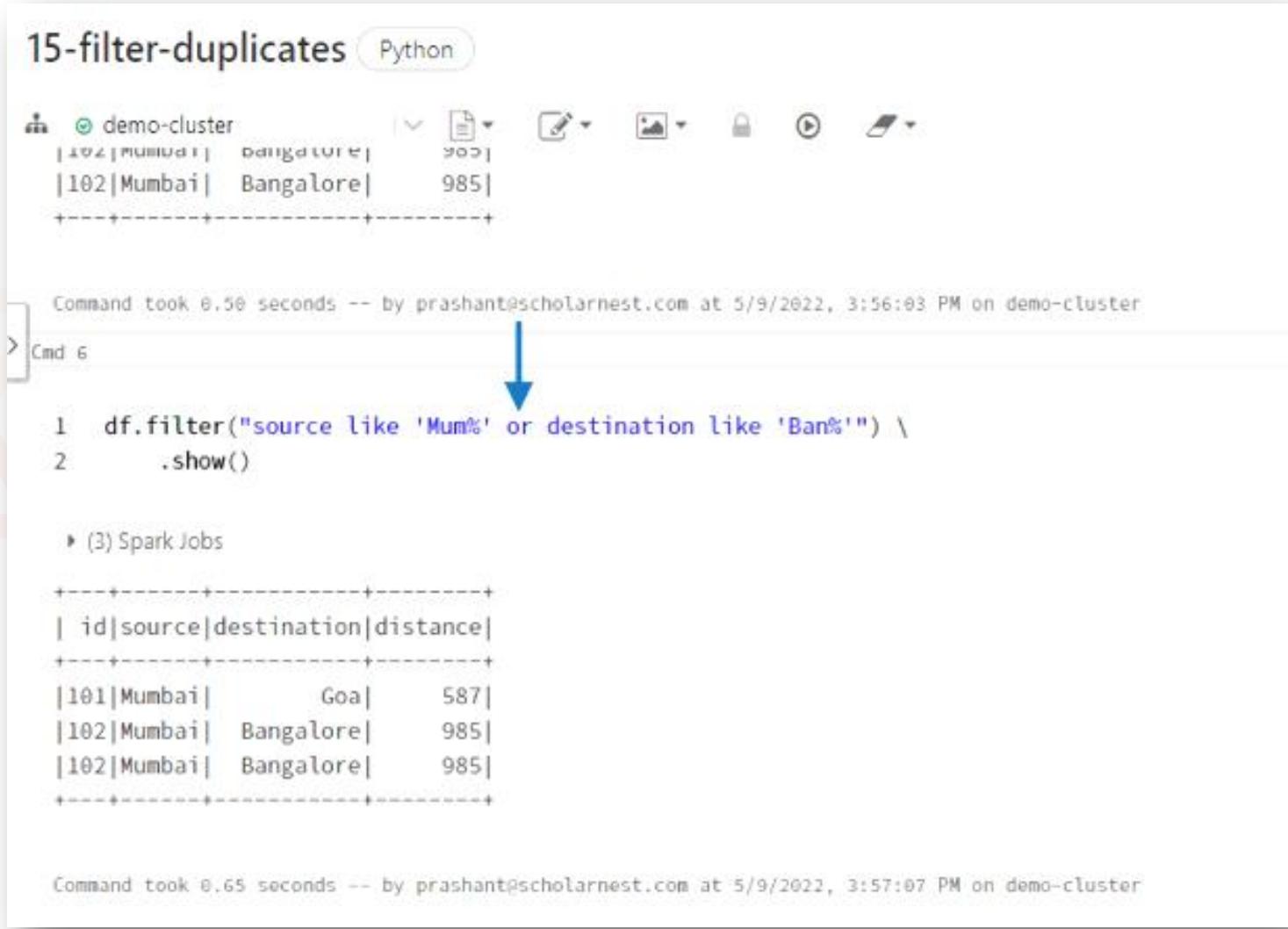
```
> Cmd 5

1 from pyspark.sql.functions import col
2
3 df.filter(col("source").like("Mum%") & col("destination").like("Ban%")) \
4 .show()

▶ (3) Spark Jobs
+---+-----+-----+
| id|source|destination|distance|
+---+-----+-----+
| 102|Mumbai| Bangalore|    985|
| 102|Mumbai| Bangalore|    985|
+---+-----+-----+

Command took 0.50 seconds -- by prashant@scholarnest.com at 5/9/2022, 3:56:03 PM on demo-cluster
```

You can also use OR condition as shown below.



The screenshot shows a Jupyter Notebook cell titled "15-filter-duplicates" in Python mode. The cell contains the following code:

```
1 df.filter("source like 'Mum%' or destination like 'Ban%'") \
2     .show()
```

A blue arrow points to the line of code where the OR condition is used. The output of the command is displayed below:

id	source	destination	distance
101	Mumbai	Goa	587
102	Mumbai	Bangalore	985
102	Mumbai	Bangalore	985

The command took 0.65 seconds to execute.

We can change the ampersand to a pipe for using the OR condition in column API.

15-filter-duplicates Python

demo-cluster

	source	destination	distance
101	Mumbai	Goa	587
102	Mumbai	Bangalore	985
102	Mumbai	Bangalore	985

Command took 0.65 seconds -- by prashant@scholarnest.com at 5/9/2022, 3:57:07 PM on demo-cluster

Cmd 7

```
1 from pyspark.sql.functions import col
2
3 df.filter(col("source").like("Mum%") | col("destination").like("Ban%")) \
4 .show()
```

▶ (3) Spark Jobs

	source	destination	distance
101	Mumbai	Goa	587
102	Mumbai	Bangalore	985
102	Mumbai	Bangalore	985

Command took 0.52 seconds -- by prashant@scholarnest.com at 5/9/2022, 3:57:44 PM on demo-cluster

Here is the Dataframe shown below. And we have few problems for duplication in this Dataframe. Let's try it one by one.

```
> [Cmd 8]
1 df.show()

▶ (3) Spark Jobs
+---+-----+-----+-----+
| id| source|destination|distance|
+---+-----+-----+-----+
|101| Mumbai|        Goa|      587|
|102| Mumbai|    Bangalore|     985|
|102| Mumbai|    Bangalore|     985|
|103| Delhi|     Chennai|    2208|
|104| Delhi|     Chennai|    2208|
|105| Bangalore|     Kolkata|   1868|
|105| Bangalore|     Kolkata|   1865|
+---+-----+-----+-----+
Command took 0.65 seconds -- by prashant@scholarnest.com at 5/11/2022, 9:26:10 PM on demo-cluster
[Cmd 9]
```

102 - Entire row is duplicate
105 - Record ID is duplicate
103/104 - ID is different, but other fields are duplicate

Here is the solution to the first problem. So one record is gone.

You can see only one record for 102.

The distinct() method is similar to the DISTINCT in SQL queries.

If your record is duplicated across all the columns, you can use the distinct() method.

```
> Cmd 9  
1 df.distinct().show()  
  
▶ (2) Spark Jobs  
  
+---+-----+-----+-----+  
| id| source|destination|distance|  
+---+-----+-----+-----+  
| 101| Mumbai|        Goa|      587|  
| 102| Mumbai|    Bangalore|     985|  
| 103|   Delhi|      Chennai|    2208|  
| 104|   Delhi|      Chennai|    2208|  
| 105|Bangalore|      Kolkata|   1868|  
| 105|Bangalore|      Kolkata|   1865|  
+---+-----+-----+-----+
```

Command took 1.10 seconds -- by prashant@scholarnest.com at 5/11/2022, 9:26:42 PM on demo-cluster

We also have one more alternative for the first problem, the dropDuplicates() method. So the result of the dropDuplicates() and the distinct() is the same.

```
> [Cmd 10]

1 df.dropDuplicates().show()

▶ (2) Spark Jobs

+---+-----+-----+-----+
| id| source|destination|distance|
+---+-----+-----+-----+
|101| Mumbai|      Goa|     587|
|102| Mumbai|   Bangalore|    985|
|103|   Delhi|    Chennai|   2208|
|104|   Delhi|    Chennai|   2208|
|105|Bangalore|    Kolkata|   1868|
|105|Bangalore|    Kolkata|   1865|
+---+-----+-----+-----+

Command took 0.68 seconds -- by prashant@scholarnest.com at 5/11/2022, 9:27:28 PM on demo-cluster
```

Now, look at the second problem. The record id 105 is duplicate, but other column values are different. This one is not a whole record duplication. So the distinct() method does not help. However, the dropDuplicates() method is more powerful. It can take a list of columns and eliminate duplicates based on those columns as shown below. So the drop dropDuplicates() will look at the id column and remove all records where the id is duplicate. And this guy solves two problems in one shot. It eliminated a duplicate of 102 and also 105. The 102 was a complete duplicate. That one is gone because the id was also duplicated. The 105 is fine because the id was duplicated.

```
> Cmd 11

1 df.dropDuplicates(["id"]).show()

▶ (2) Spark Jobs

+---+-----+-----+
| id | source|destination|distance|
+---+-----+-----+
| 101 | Mumbai|      Goa|     587|
| 102 | Mumbai|   Bangalore|    985|
| 103 | Delhi|     Chennai|   2208|
| 104 | Delhi|     Chennai|   2208|
| 105 | Bangalore|   Kolkata|  1868|
+---+-----+-----+

Command took 1.31 seconds -- by prashant@scholarnest.com at 5/11/2022, 9:28:03 PM on demo-cluster

Cmd 12
```

Now let's go to the third problem.

103 and 104 are duplicates, but the ID is different.

How do we eliminate it? We can use the dropDuplicates() once again. But this time, we want to remove duplicates on all columns except the id column. Let me add another dropDuplicates() as shown below.

```
> Cmd 11

1 df.dropDuplicates(["id"]) \
2   .dropDuplicates(["source", "destination", "distance"]) \
3   .show()

▶ (3) Spark Jobs

+---+-----+-----+-----+
| id| source|destination|distance|
+---+-----+-----+-----+
| 102| Mumbai| Bangalore|    985|
| 105|Bangalore|      Kolkata|   1868|
| 101| Mumbai|        Goa|     587|
| 103|   Delhi|    Chennai|   2208|
+---+-----+-----+-----+

Command took 2.21 seconds -- by prashant@scholarnest.com at 5/11/2022, 9:28:56 PM on demo-cluster
```

So we learned two methods to remove duplicates.

1. `distinct()`
2. `dropDuplicates()`

The `distinct` method is similar to `DISTINCT` in the SQL, and it does not support complex scenarios.

The `dropDuplicates()` is more powerful. It takes a list of columns, compares records only using the given columns, and eliminates duplicates.



Thank You
ScholarNest Technologies Pvt Ltd.
www.scholarnest.com