

# Big Data Tools and Techniques

Chapter 7: Introduction to Spark (Part C)

# Introduction to Databricks

---

## Integrated Workspace

Notebooks

Dashboards

## BI Tools



## Your Custom Spark Apps

Production Jobs

Orchestrated Apache® Spark™ in the Cloud

Open Source  +  **databricks**™ Managed Services

## Your Storage



Cloud Storage | Data Warehouses | Data Lakes



Important Notice: [Acceptable use and unused account termination policy](#) and [Terms of Use](#) update.



### Sign In to Databricks Community Edition



Email / Username



Password

[Forgot Password?](#)

[Sign In](#)

New to Databricks? [Sign Up.](#)

[Privacy Policy](#) | [Terms of Use](#)



Search for anything



ENG

2:55 PM  
6/6/2020



Login - Databricks Community

Try Databricks

databricks.com/try-databricks

Platform Solutions Customers Learn Partners Events Open Source Company

EN SUPPORT CONTACT LOG IN TRY DATABRICKS

# Try Databricks

An open and unified data analytics platform for data engineering, machine learning, and analytics

From the original creators of Apache Spark™, Delta Lake, MLflow, and Koalas

Tell us a little about yourself to get started.

\* First Name:

\* Last Name:

\* Company Name

\* Work Email

\* How would you describe your role?

Select...

\* What is your intended use case?

Select...

Quickstart Notebook.sql

Show all



Search for anything



2:55 PM

6/6/2020

1

# Welcome to databricks

[Upgrade](#)

## Explore the Quickstart Tutorial

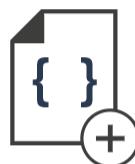
Spin up a cluster, run queries on preloaded data, and display results in 5 minutes.

Drop files or [click to browse](#)



## Import & Explore Data

Quickly import data, preview its schema, create a table, and query it in a notebook.



## Create a Blank Notebook

Create a notebook to start querying, visualizing, and modeling your data.

## Common Tasks

- New Notebook
- Create Table
- New Cluster
- New Job
- New MLflow Experiment
- Import Library
- Read Documentation

## Recents

- Quickstart Notebook

## What's new in v3.21

[View latest release notes](#)



## Workspace

Users

Quickstart Notebook

atabricks

Upgrade

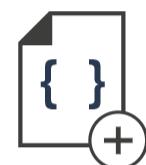
?

User icon

Drop files or [click to browse](#)

## Import &amp; Explore Data

Import data from various sources and explore it using Databricks' built-in tools.



## Create a Blank Notebook

Create a new notebook to start querying, visualizing, and modeling your data.

## Recents

Quickstart Notebook

## What's new in v3.21

[View latest release notes](#)

Search for anything

3:11 PM  
6/6/2020

1



## Recent

[Quickstart Notebook](#)[Upgrade](#)Welcome to **databricks**[Explore the Quickstart Tutorial](#)

Enter, run queries on preloaded data, and display results in 5 minutes.

[Import & Explore Data](#)

Quickly import data, preview its schema, create a table, and query it in a notebook.

[Create a Blank Notebook](#)

Create a notebook to start querying, visualizing, and modeling your data.

KS

Recents

What's new in v3.21

[Quickstart Notebook](#)[View latest release notes](#)

Experiment

try

mentation



Search for anything



3:12 PM

6/6/2020





Search Workspace

Upgrade

?

User icon

# Welcome to databricks



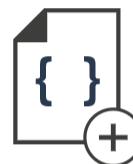
## Explore the Quickstart Tutorial

Enter, run queries on preloaded data, and display results in 5 minutes.



## Import & Explore Data

Quickly import data, preview its schema, create a table, and query it in a notebook.



## Create a Blank Notebook

Create a notebook to start querying, visualizing, and modeling your data.

KS

Recents

What's new in v3.21

Quickstart Notebook

[View latest release notes](#)

My First Notebook



Quickstart Notebook.sql



Show all



3:12 PM

6/6/2020





## Create Cluster

## New Cluster

Cancel

Create Cluster

0 Workers: 0.0 GB Memory, 0 Cores, 0 DBU  
1 Driver: 15.3 GB Memory, 2 Cores, 1 DBU 

Cluster Name

My First Spark Cluster

UI | JSON

Databricks Runtime Version

Runtime: 6.5 (Scala 2.11, Spark 2.4.5)

New This Runtime version supports only Python 3.

Instance

Free 15GB Memory: As a Community Edition user, your cluster will automatically terminate after an idle period of two hours.  
For more configuration options, please [upgrade your Databricks subscription](#).

Instances Spark

Availability Zone

us-west-2c



Search for anything



ENG

3:13 PM  
6/6/2020

Clusters - Databricks Community X +

community.cloud.databricks.com/

Clusters

+ Create Cluster

All Created by me Filter

1 clusters, 0 pinned

All-Purpose Clusters

Name	State	Nodes	Driver	Worker	Runtime	Creator	Actions
My First Spark Cluster	Pending ⓘ	0	Community Optimizer	Community Optimizer	6.5 (includes Apache ..)	0	...

Job Clusters

No clusters found

databricks

Home

Workspace

Recents

Data

Clusters

Jobs

Search



## Workspace

Users

Mode ▾ Permissions Run All Clear ▾

Publish Comments Runs Revision history

Quickstart Notebook

Create

- Clone
  - Import
  - Export
- Permissions

Notebook

Library

Folder

MLflow Experiment

e link in a new window.

(Scala 2.11, Spark 2.4.4).

**er and run all commands in the notebook**

art.

**ole from a Databricks dataset**



## Workspace

Users

Quickstart Notebook

Mode ▾ Permissions Run All Clear ▾

Publish Comments Runs Revision history

## Create Notebook

Name Default Language Cluster 

Cancel

Create

er and run all commands in the notebook

art.

ole from a Databricks dataset

Quickstart Notebook.sql

Show all

X



Search for anything



3:22 PM

6/6/2020



## Create Library

## Library Source

Upload DBFS/S3 PyPI Maven CRAN

## Package

PyPI package (simplejson or simplejson==3.8.0)

## Repository ?

Optional

CreateCancel



## My First Notebook (Python)

█ My First Spark Cluster █ File █ Edit █ View: Code █ Permissions █ Run All █ Clear █ Publish █ Comments █ Runs █ Revision history

Cmd 1  

```
1
```

▶ ▷ ⏪ ⏴ ⏵ ⏳

Shift+Enter to run    shortcuts

New Notebook

Clone

Rename

Move

Delete

**Export** ➔

Publish

Clear Revision History

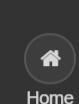
Change Default Language

DBC Archive

Source File

**IPython Notebook**

HTML



Quickstart Notebook.sql



Show all



Search for anything

3:23 PM  
6/6/2020



## Workspace

Users



Home



Workspace



Recents



Data



Clusters



Jobs



Search

## Import Notebooks

Import from:  File  URLDrop files to upload, or [browse](#).

Accepted formats: .dbc, .scala, .py, .sql, .r, .ipynb, .Rmd, .html

(To import a library, such as a jar or egg, [click here](#))

Cancel

Import

**er and run all commands in the notebook**

art.

**ole from a Databricks dataset**

Search for anything





Detached

File ▾

Edit ▾

View: Code ▾

Permissions

Run All

Clear ▾



Publish

Comments

Runs

Revision history

Cmd 1

## Databricks in 5 minutes



Cmd 2

### Create a quickstart cluster

1. In the sidebar, right-click the **Clusters** button and open the link in a new window.
2. On the Clusters page, click **Create Cluster**.
3. Name the cluster **Quickstart**.
4. In the Databricks Runtime Version drop-down, select **6.3 (Scala 2.11, Spark 2.4.4)**.
5. Click **Create Cluster**.

Cmd 3

### Attach the notebook to the cluster and run all commands in the notebook

1. Return to this notebook.
2. In the notebook menu bar, select **Detached ▾ > Quickstart**.
3. When the cluster changes from  to , click **Run All**.

Cmd 4

### The next command creates a table from a Databricks dataset

Cmd 5

```
1 DROP TABLE IF EXISTS diamonds;
2
3 CREATE TABLE diamonds
```



Search for anything



ENG

3:05 PM  
6/6/2020

Quickstart Notebook - Databricks X +

community.cloud.databricks.com/

Quickstart Notebook (SQL)

Detached File Edit View: Code Permissions Run All Clear Publish Comments Runs Revision history

Cmd 4

The next command creates a table from a Databricks dataset

Cmd 5

```
1 DROP TABLE IF EXISTS diamonds;
2
3 CREATE TABLE diamonds
4 USING CSV
5 OPTIONS (path "/databricks-datasets/Rdatasets/data-001/csv/ggplot2/diamonds.csv", header "true")
6
```

OK

Command took 46.18 seconds -- by a user at 6/13/2019, 5:06:53 AM on unknown cluster

Cmd 6

```
1 SELECT * from diamonds
```

Showing the first 1000 rows.

_c0	carat	cut	color	clarity	depth	table	price	x	y	z
1	0.23	Ideal	E	SI2	61.5	55	326	3.95	3.98	2.43
2	0.21	Premium	E	SI1	59.8	61	326	3.89	3.84	2.31
3	0.23	Good	E	VS1	56.9	65	327	4.05	4.07	2.31
4	0.29	Premium	I	VS2	62.4	58	334	4.2	4.23	2.63
5	0.24	Good	E	SI1	63.3	56	325	4.01	4.05	2.75

Command took 6.11 seconds -- by a user at 6/13/2019, 5:06:53 AM on unknown cluster

Cmd 7

```
1 %python
```

Search for anything

File Explorer Mail Word Chrome Powerpoint Edge Task View Home Workspace Recents Data Clusters Jobs Search

3:06 PM  
6/6/2020

Quickstart Notebook - Databricks X +

community.cloud.databricks.com/

Quickstart Notebook (SQL)

Detached File Edit View: Code Permissions Run All Clear Publish Comments Runs Revision history

Cmd 7

```
1 %python
2 diamonds = spark.read.csv("/databricks-datasets/Rdatasets/data-001/csv/ggplot2/diamonds.csv", header="true", inferSchema="true")
3 diamonds.write.format("delta").save("/delta/diamonds")
```

▶ diamonds: pyspark.sql.dataframe.DataFrame = [c0: integer, carat: double ... 9 more fields]

Command took 13.50 seconds -- by a user at 6/13/2019, 5:06:53 AM on unknown cluster

Cmd 8

```
1 DROP TABLE IF EXISTS diamonds;
2
3 CREATE TABLE diamonds USING DELTA LOCATION '/delta/diamonds/'
```

OK

Command took 0.77 seconds -- by a user at 6/13/2019, 5:06:53 AM on unknown cluster

Cmd 9

```
1 SELECT * from diamonds
```

_c0	carat	cut	color	clarity	depth	table	price	x	y	z
1	0.23	Ideal	E	SI2	61.5	55	326	3.95	3.98	2.43
2	0.21	Premium	E	SI1	59.8	61	326	3.89	3.84	2.31
3	0.23	Good	E	VS1	56.9	65	327	4.05	4.07	2.31
4	0.29	Premium	I	VS2	62.4	58	334	4.2	4.23	2.63

Showing the first 1000 rows.

Command took 6.89 seconds -- by a user at 6/13/2019, 5:06:53 AM on unknown cluster

Cmd 10

Search for anything

Windows Start E Microsoft Edge File Explorer Mail Word Powerpoint Google Chrome Task View Search Help ENG 3:06 PM 6/6/2020

Quickstart Notebook - Databricks X +

community.cloud.databricks.com/

Quickstart Notebook (SQL)

Detached File Edit View: Code Permissions Run All Clear Publish Comments Runs Revision history

Command took 6.89 seconds -- by a user at 6/13/2019, 5:06:53 AM on unknown cluster

Cmd 10

## The next command manipulates the data and displays the results

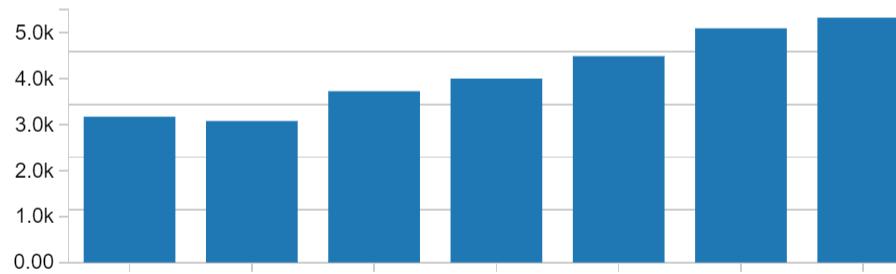
Specifically, the command:

1. Selects color and price columns, averages the price, and groups and orders by color.
2. Displays a table of the results.

Cmd 11

```
1 | SELECT color, avg(price) AS price FROM diamonds GROUP BY color ORDER BY color
```

price



color	price
D	~3.1k
E	~3.1k
F	~3.7k
G	~4.0k
H	~4.5k
I	~4.9k
J	~5.2k

Plot Options... Download

Command took 2.83 seconds -- by a user at 6/13/2019, 5:06:53 AM on unknown cluster

Cmd 12

## Convert the table to a chart

Search for anything

File Edit View: Code Permissions Run All Clear Publish Comments Runs Revision history

3:07 PM 6/6/2020 ENG 1

Quickstart Notebook - Databricks X +

community.cloud.databricks.com/ 067445025/command/1158110067445026

Detached File Edit View: Code Permissions Run All Clear Publish Comments Runs Revision history

Home Workspace Recents Data Clusters Jobs Search

## Quickstart Notebook (SQL)

### Convert the table to a chart

Under the table, click the bar chart icon.

Cmd 13

### Repeat the same operations using Python DataFrame API.

This is a SQL notebook; by default command statements are passed to a SQL interpreter. To pass command statements to a Python interpreter, include the `%python` magic command.

Cmd 14

### The next command creates a DataFrame from a Databricks dataset

Cmd 15

```
1 %python
2 diamonds = spark.read.csv("/databricks-datasets/Rdatasets/data-001/csv/ggplot2/diamonds.csv", header="true", inferSchema="true")
```

▶ diamonds: pyspark.sql.dataframe.DataFrame = [c0: integer, carat: double ... 9 more fields]

Command took 1.67 seconds -- by a user at 6/13/2019, 5:06:53 AM on unknown cluster

Cmd 16

### The next command manipulates the data and displays the results

Cmd 17

```
1 %python
2 from pyspark.sql.functions import avg
3
4 display(diamonds.select("color", "price").groupBy("color").agg(avg("price")).sort("color"))
```

Search for anything

Windows Start button

Icons: File Explorer, Microsoft Edge, Microsoft Store, Mail, Word, OneDrive, Google Chrome, Microsoft Edge Dev, Task View, Question mark, Volume, Network, Cloud, Task Manager, ENG, 3:08 PM, 6/6/2020, 1

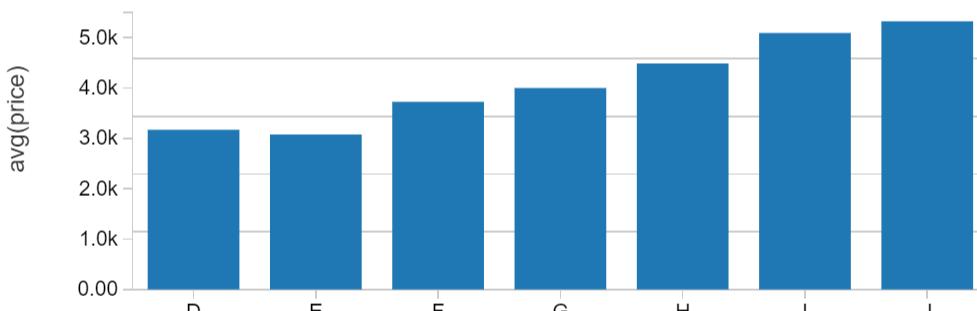
## Quickstart Notebook (SQL)

Detached File Edit View: Code Permissions Run All Clear Publish Comments Runs Revision history

### The next command manipulates the data and displays the results

Cmd 17

```
1 %python
2 from pyspark.sql.functions import avg
3
4 display(diamonds.select("color","price").groupBy("color").agg(avg("price")).sort("color"))
```



Plot Options... Download

Command took 2.24 seconds -- by a user at 6/13/2019, 5:06:53 AM on unknown cluster

Shift+Enter to run [shortcuts](#)



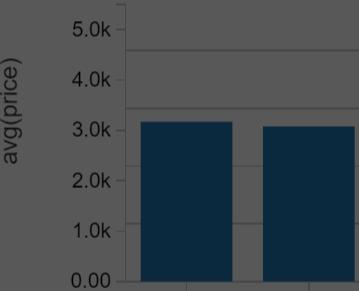
## Quickstart Notebook (SQL)

Detached

### The next command

Cmd 17

```
1 %python
2 from pyspark.sql.functions
3
4 display(diamonds.select("co
```



Plot Options...

Command took 2.24 seconds -- by a user at 6/13/2019, 5:06:53 AM on unknown cluster

Shift+Enter to run [shortcuts](#)

Quickstart Notebook.sql

Show all

X



Search for anything



3:25 PM  
6/6/2020



### Customize Plot

All fields:

color  
avg(price)  
<id>

Keys:

color x

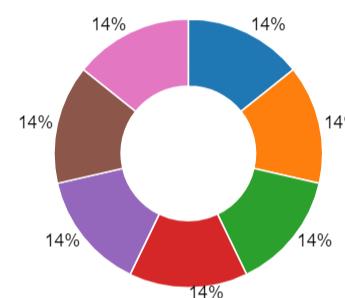
Series groupings:

Values:

avg(price) x

color

- D
- E
- F
- G
- H
- I
- J



Donut

Aggregation: COUNT

Display type: Pie chart

Cancel

Apply

Attached: reporting

View: Code

File

Permissions

Run All



Comments

Revision history



databricks



Home



Workspace



Recent



Tables



Clusters



Jobs



Apps



Search



Settings

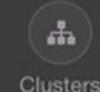
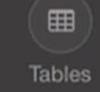
```
> select m.ClientID, c.CountryCode3, m.SessionID  
  from mobile_sample m  
  join countrycodes c  
    on m.Country = c.CountryName
```



## ▶ (4) Spark Jobs

ClientID	CountryCode3	SessionID	DeviceMake
96046	ARG	0	ASUS
96987	ARG	0	ASUS
127767	ARG	0	ASUS
3846	AUS	0	SAMSUNG
8509	AUS	0	Apple
10085	AUS	0	HTC
14719	AUS	0	HTC
15286	AUS	0	Apple
15286	AUS	0	Apple
15286	AUS	0	Apple
15286	Bar	0	Apple
15286	Scatter	0	Apple
15286	Map	0	Apple
15286	Line	0	Apple
15286	Area	0	Apple
15286	Pie	0	Apple

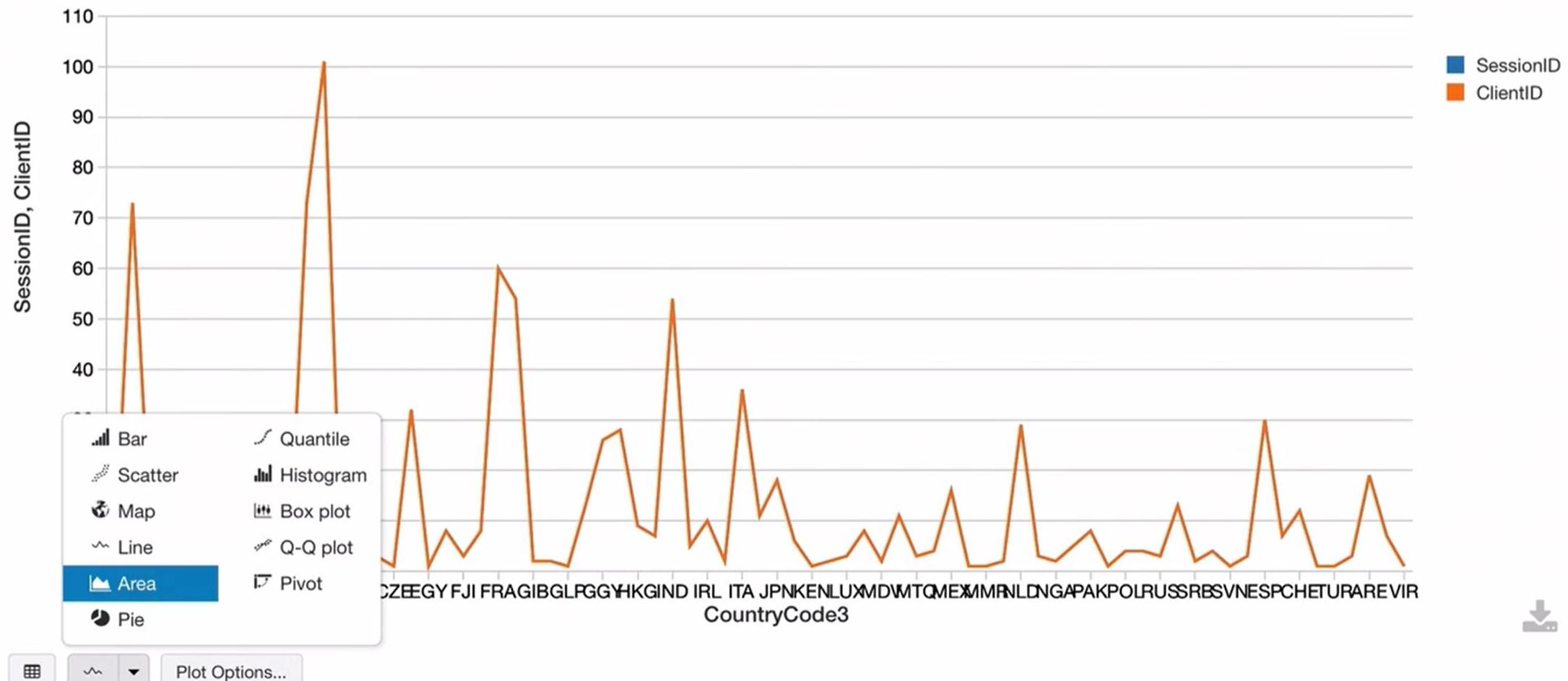




```
> select m.ClientID, c.CountryCode3, m.SessionID  
from mobile_sample m  
join countrycodes c  
on m.Country = c.CountryName
```



## ▶ (4) Spark Jobs





databricks



Home



Workspace



Recent



Tables



Clusters



Jobs



Apps



Search



Settings

Attached: reporting

View: Code

File

Permissions

Run All



Comments

Revision history

```
> select m.ClientID, c.CountryCode3, m.SessionID  
from mobile_sample m  
join countrycodes c  
on m.Country = c.CountryName
```

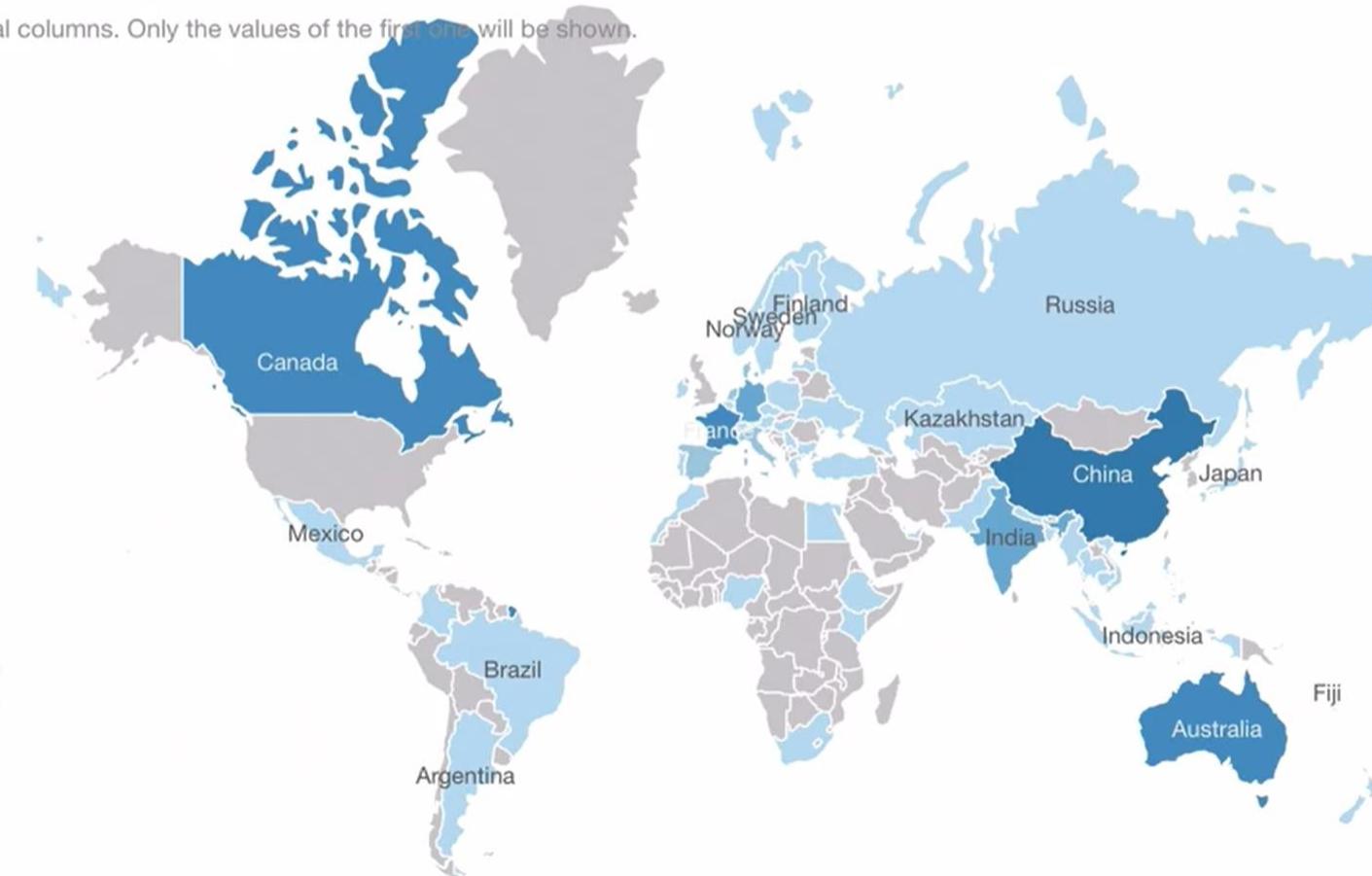


-

x

## ▶ (4) Spark Jobs

More than two numerical columns. Only the values of the first one will be shown.



Plot Options...