# Java-Success.com

Prepare to fast-track, choose & go places with 800+ Java & Big Data Q&As with lots of code & diagrams.

| search here … | Go |

Home    Why? ▼    300+ Java FAQs ▼    300+ Big Data FAQs ▼    Courses ▼

👤 Membership ▼    Your Career ▼

01: Databricks getting started – Spark, Shell, SQL

# 01: Databricks getting started – Spark, Shell, SQL

📅 Posted on April 10, 2020

**Step 1:** Signup to Databricks community edition – https://databricks.com/try-databricks. Fill in the details and you can leave your mobile number blank. Select "COMMUNITY EDITION" ==> "GET STARTED".

## 300+ Java Interview FAQs

300+ Java FAQs 🔥    ⌄

16+ Java Key Areas Q&As    ⌄

150+ Java Architect FAQs    ⌄

80+ Java Code Quality Q&As    ⌄

150+ Java Coding Q&As    ⌄

## 300+ Big Data Interview FAQs

300+ Big Data FAQs 🔥    ⌄

Tutorials - Big Data    ›

TUT - 🔲 Starting Big Data

TUT - Starting Spark & Scala

Databricks getting started

If you have a Cloud account then you can use it.

**Step 2:** Check your email and click the "link" in the email & reset your password.

**Step 3:** Login to Databricks notebook:

https://community.cloud.databricks.com/login.html.



Databricks getting started dashboard

**Step 4:** Create a **CLUSTER** and it will take a few minutes to come up. This cluster will go down after 2 hours.

## 800+ Java Interview Q&As

300+ Core Java Q&As ⌄

300+ Enterprise Java Q&As ⌄

150+ Java Frameworks Q&As ⌄

120+ Companion Tech Q&As ⌄

Tutorials - Enterprise Java ⌄

Databricks create a cluster

## Step 5: Select "DATA", and upload a file named "employee.csv".

```
1  emp_id,emp_name,emp_city,emp_salary
2  1, John,  Sydney,  35000.00
3  2, Peter, Melbourne,  45000.00
4  3, Sam,  Sydney,55000.00
5
```



Databricks upload a .csv

**Step 6:** "Create Table With UI" as shown below:
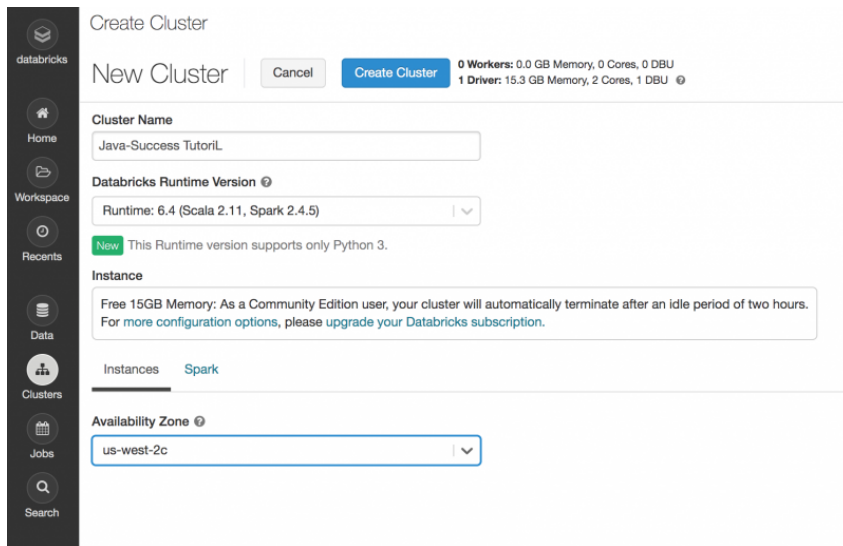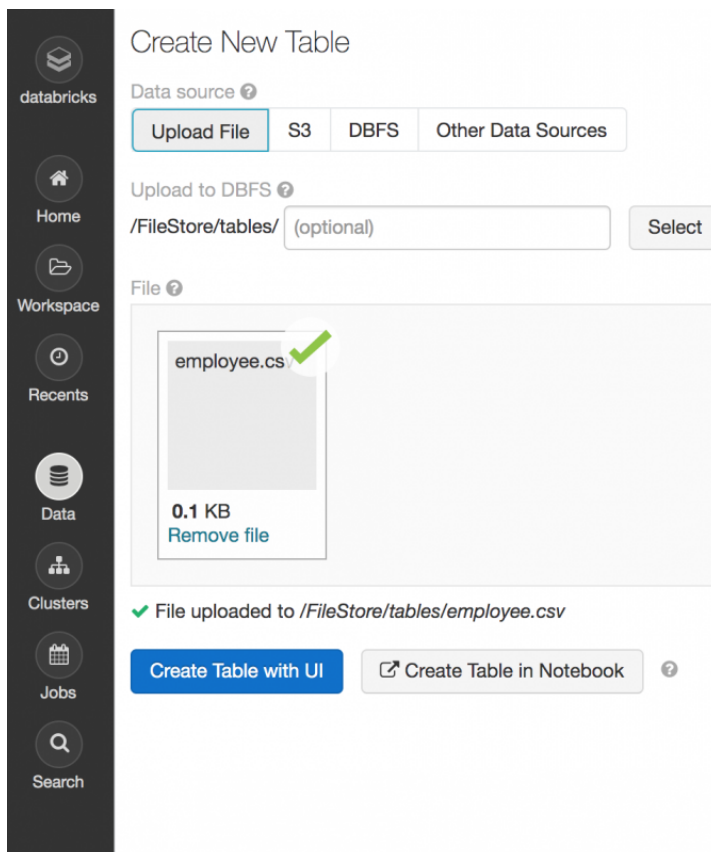
**Note:** Please check the "First row is header" check box on the LHS so that column names appear from the file.



Databricks create a table with .csv

Click on "Create Table".

**Step 7:** Click on the "databricks" icon on the LHS menu, and then "Create a Blank Notebook".



Databricks blank notebook

# Spark in Python (i.e.PySpark)

Since we created the notebook as "**python**", we don't
have to do "**%python**" as it is the default language. If
you want to use "scala" then add "**%scala**" as the first
line in a **cell**.

**Step 8:** Run the below PySpark code to display
uploaded "**/FileStore/tables/employee.csv**" as a
**Dataframe**.

```
1  employees_df = spark.read.csv("/FileStore/tables/er
2  employees_df.show()
3
```

Click on "**Run**" in the top "**RHS**" menu.



Databricks Notebook simple Dataframe

▶▼    ✔    ━    ✖

⧉ Copy Cell

✂ Cut Cell

⬇ Export Cell

⬛ Add Cell Above

⬛ Add Cell Below

↓ Move Down

H Show Title

</> Hide Code

⟨/⟩ Hide Result

Use "Down Arrow" on in the RHS of a cell to create a new cell.

**Step 9:** Add a new column to the Dataframe in a separate cell and then **Run**.

```
1  from pyspark.sql import functions as F
2
3  employees_df.withColumn("updated_ts", F.current_tir
4
```

Databricks notebook add new column to dataframe

# Run a shell command

You can run a shell command with "%sh" as shown below:

```
1  %sh
2
3  ls -ltr
4
```

You can click on the "**+**" in the middle to add new cells.



Databricks shell command

# Run a SQL command

```
1  %sql
2
3  SELECT * FROM employee_csv;
4
```

Databricks SQL command

# dbutils

```python
%python

dbutils.fs.rm("/tmp/my-delta-table", recurse=True)
```

# Spark in Scala

```scala
%scala

val df_employees = sqlContext.read.format("com.dat
    .option("delimiter", ",")
    .load("/FileStore/tables/employee.csv")

import org.apache.spark.sql.functions._
df_employees.withColumn("updated_ts", lit(current_
```



Databricks Spark in Scala

# Important: PySpark API

Have the PySpark API PySpark modules handy to
code. You can click on "Dataframe" to see what

functions are available. For example, withColumn function in the Dataframe module



**withColumn**(colName, col)                                                          [source]

Returns a new `DataFrame` by adding a column or replacing the existing column that has the same name.

The column expression must be an expression over this `DataFrame`; attempting to add a column from some other `DataFrame` will raise an error.

Parameters:   • **colName** – string, name of the new column.
              • **col** – a `Column` expression for the new column.

```
>>> df.withColumn('age2', df.age + 2).collect()
[Row(age=2, name='Alice', age2=4), Row(age=5, name='Bob', age2=7)]
```

*New in version 1.3.*

Databricks PySpark Modules & API

## Where are my notebooks saved?

Your notebooks will be saved under "**workspace**" ==> "**users**" ==> "**[your username]**"

## What if you want to practice in Scala?

You an try the examples from Tutorials – Spark Scala on Zeppelin with some minor changes.

‹   AWS Web Application Security Q&As

           02: Databricks – Spark schemas, casting & PySpark API   ›

## Disclaimer

The contents in this Java-Success are copyrighted and from EmpoweringTech pty ltd. The EmpoweringTech pty ltd has the right to correct or enhance the current content without any prior notice. These are general advice only, and one needs to take his/her own circumstances into consideration. The EmpoweringTech pty ltd will not be held liable for any damages caused or alleged to be caused either directly or indirectly by these materials and resources. Any trademarked names or labels used in this blog remain the property of their respective trademark owners. Links to external sites do not imply endorsement of the linked-to sites. Privacy Policy

© 2022 java-success.com