800+ Q&As | Logout | Contact

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 ▼

Home > bigdata-success.com > Tutorials - Big Data > TUT - Spark Scala on Zeppelin >

05: Spark on Zeppelin - semi-structured log file

05: Spark on Zeppelin – semi-structured log file



This tutorial extends the series: Spark on Apache Zeppelin Tutorials.

Step 1: Pull apache/zeppelin image from the docker hub, and build the image with the following command.

```
1 $ docker pull apache/zeppelin:0.7.3
```

"docker images" will show the image that was created.

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

Step 2: Run the above image to create a container with the following command.

```
1 $ docker run -it -p 8080:8080 apache/zeppelin:0.7.3
```

Step 3: Go to a browser and type:

"http://locahost:8080" to verify if zeppelin note book is accessible.

Step 4: Create a ssh session into the docker container from a terminal:

Firstly, get the container id with:

```
1 $ docker ps
2
```

Use the container id E.g. 89bce1cdc799 to ssh.

```
1 $ docker exec -it 89bce1cdc799 /bin/bash 2
```

Step 5: Once inside the docker container running on ubuntu, install "vim" editor.

```
1 | $ apt-get update
2 | $ apt-get install vim
3 |
```

Step 6: Create a temp.log file with vim.

```
1 $ vim temp.log
2
```

TUT - Starting with Python

TUT - Kafka

TUT - Pig

TUT - Apache Storm

TUT - Spark Scala on Zeppelin

TUT - Cloudera

TUT - Cloudera on Docker

TUT - File Formats

TUT - Spark on Docker

TUT - Flume

TUT - Hadoop (HDFS)

TUT - HBase (NoSQL)

TUT - Hive (SQL)

TUT - Hadoop & Spark

TUT - MapReduce

TUT - Spark and Scala

TUT - Spark & Java

TUT - PySpark on Databricks TUT - Zookeeper

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



Enter the following text by pressing "i" to insert.

```
1 ERROR 2009-09-19 09:56:01 main RDFDefaultErrorHand INFO 2009-09-13 09:56:01 login Some.java:44 Some INFO 2009-09-11 09:56:01 login Some.java:44 Some WARNING 2009-09-09 09:56:01 main MarketData.java:45 INFO 2009-09-06 09:56:01 login Some.java:44 Some INFO 2009-09-05 09:56:01 login Some.java:44 Some INFO 2009-09-02 09:56:01 login Some.java:44 Some INFO 2009-09-01 09:56:01 login Some.java:44 Some INFO 2009-09-01 09:56:01 login Some.java:44 Some ERROR 2009-08-25 09:56:01 main RDFDefaultErrorHand
```

exit vim by pressing "wq!".

Step 7: On zeppelin note book on the browser http://localhost:8080, create a note book and use "%sh" to list the files.

Step 8: Create a spark RDD by reading this temp.log file. The path to temp.log file is "file:///zeppelin/temp.log"

Step 9: Create a DataFrame as shown below:

```
1
```

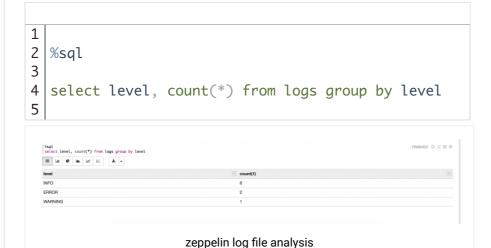
```
2
   %spark
3
4
   import java.text.SimpleDateFormat
5
   import java.sql.Date
6
7
   val DATE_FORMAT = new SimpleDateFormat("yyyy-mm-de
8
   case class Log(level: String, date: java.sql.Date
9
10
   val logsDf = logsRdd.map { line =>
11
12
13
       val s = line.split(" ")
14
       val logLevel = s(0)
       val dateTime = DATE_FORMAT.parse(s(1) + " " +
15
       val fileName = s(3).split(":")(0)
16
17
       Log(logLevel,new Date(dateTime.getTime()), fil
18
19 \ \ . toDF()
20
21 logsDf.show()
22
```

Output:

```
1
2
3
      level
                 date|fileName|
4
5
      ERROR | 2009-01-19 |
                          main
6
       INFO|2009-01-13|
                         login
7
       INFO|2009-01-11| login|
8
   |WARNING|2009-01-09|
                         main|
9
       INFO|2009-01-06|
                         login
10
       INFO|2009-01-05| login|
11
      INFO|2009-01-02| login|
12
      INFO|2009-01-01| login|
13
  | ERROR|2009-01-25|
                         main|
14
15
```

Step 10: Register as temp table so that you can run SQL queries to analyze the data.

Step 11: Analyse with SQL – tabular and graph.



04: Spark on Zeppelin – DataFrame joins in Scala40+ Apache Spark best practices & optimisation interview FAQs – Part-1

>

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