Menu | Logout

# Java-Success.com

800+ Java & Big Data Interview Q&As with code & diagrams to fast-track & go places with choices.

search here …                    Go

Home    300+ Java FAQs ▾    300+ Big Data FAQs ▾    Courses ▾    👤 Membership ▾    Career ▾

Home › bigdata-success.com › 300+ Big Data FAQs 🔥 › FAQs Data - 08: Spark › 17: Spark interview Q&As with coding examples in pyspark (i.e. python)

# 17: Spark interview Q&As with coding examples in pyspark (i.e. python)

📅 Posted on May 13, 2018

**Q01.** How will you create a Spark context?
**A01.**

```
1  from pyspark.sql import SparkSession
2
3  spark = SparkSession.builder.appName("my spark job")
4  spark.master('local[*]')
5
6  spark.config('spark.jars.packages', 'com.amazonaws:aws-java-sdk:1.11
7           .config('spark.hadoop.mapreduce.fileoutputcommitter.alg
8           .config('spark.speculation', 'false')
9
```

**Q02.** How will you create a Dataframe by reading a file from AWS S3 bucket?
**A02.**

```
1
2      csvFileAsDataframe = spark.read.format("com.databricks.spark.cs
3               .option("header", "false") \
4               .option("inferSchema", "true") \
5               .load(s3://my-bucket/some-path/input-file.cs
6
```

**Q03.** How will you create a Dataframe by reading a table in a database?
**A03.**

```
1
2      jdbcTableAsDataframe = self.spark.read \
3               .format("jdbc") \
4               .option("url", "jdbc:mysql://myhost:3306/mydatabas
5               .option("dbtable", "mytable") \
6               .option("user", "root") \
```

## 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 🔥                    ❯

FAQs Data - 01: SQL

FAQs Data - 02: Data Modelling

FAQs Data - 02: Data Warehouse
FAQs Data - 03: Big Data

FAQs Data - 04: Hadoop (HDFS)

FAQs Data - 05: MapReduce

FAQs Data - 06: Hive

FAQs Data - 07: Impala

FAQs Data - 08: Spark

FAQs Data - 09: Spark SQL

FAQs Data - 10: Apache Kafka

FAQs Data - 11: Data Governance
FAQs Data - 11: NoSQL

FAQs Data - 12: Data security

FAQs Data - 13: Analytics & Science
FAQs Data - 14: AWS

FAQs Data - 15: Sqoop & Nifi

FAQs Data - 16: Yarn, Zookeeper

FAQs Data - 40: Scala
140+ FAQs                                ⌄

```
  7                .option("password", "password") \
  8                .option("driver", "com.mysql.jdbc.Driver") \
  9                .option("useSSL", false) \
 10                .load()
 11
```

Tutorials - Big Data

**Q04.** How will you select columns from a Dataframe?
**A04.**

```
  1
  2      from pyspark.sql.functions import col
  3
  4      outputDataframe = csvFileAsDataframe.select(col("_c0").alias("em
  5                                          col("_c1").alias("em
  6                                          col("_c2").alias("de
  7                                          col("_c3").alias("sa
  8
```

## 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

**Q05.** How do you join two Dataframes and filter?
**A05.**

```
  1
  2  #join condition
  3  condSendDate = [dfTable1.party_id == dfTable2.send_party_id, dfTable
  4
  5  #joining 3 tables
  6  dfJoined = dfTable1.join(dfTable3, dfTable1.campaign_id== dfTable3.l
  7              .join(dfAction, dfTable1.action_id == dfAction.action
  8              .join(dfTable2, condSendDate , 'left_outer' ) \
  9              .filter(dfTable3['channel'].rlike("[wW][eE][bB]")) \
 10              .filter(to_timestamp(dfTable1.available_from_date) <=
 11              .filter(to_timestamp(dfTable1.available_until_date) >=
 12
```

**Q06.** How will you add a new column to a Dataframe?
**A06.** Use "withColumn" function.

```
  1
  2  from pyspark.sql.functions import dense_rank
  3  from pyspark.sql.window import Window
  4
  5  #rank is a new column evaluated
  6  dfTable1.withColumn("rank", dense_rank().over(windowSpec))
  7
```

**Q07.** What is a window function in Spark?
**A07.** A **window function** calculates a return value for every input row of a table based on a group of rows. Spark SQL supports three kinds of window functions: ranking functions, analytic functions, and aggregate functions.

Students are grouped by student_ids and for each student his/her subjects are ranked by "written_test_grade" and "practical_test_grade" in descending order and thetop ranking subject is selected.

```
  1
  2  from pyspark.sql.window import Window
  3  from pyspark.sql.functions import dense_rank
  4
  5  windowSpec = Window.partitionBy(dfTable1['student_id']) \
```

```
 6              .orderBy(dfTable1['written_test_grade'].desc(), dfTable1[
 7
 8  dfTable1Ranked = dfTable1.withColumn("rank", dense_rank().over(windo
 9  dfTable1Toprank =  dfTable1Ranked.where(dfTable1Ranked.rank == 1)
10
```

**Q08.** When performing Dataframe operations, how do you verify your results?
**A08.** The "show()" and "printSchema()" methods are very handy for debugging.

```
 1
 2  myDataframe.show()      # prints the top 20 rows in the dataframe
 3  myDataframe.show(100)   # prints the top 100 rows in the dataframe
 4
 5  myDataframe.printSchema()   # prints the schema of the Dataframe
 6
```

**Q09.** What is a Spark udf?
**A09.** udf stands for User Defined Functions. Here is an example to create a new column by adding a few days to a given date column.

Spark let's you define custom SQL functions called user defined functions (UDFs). UDFs are great when built-in SQL functions aren't sufficient, but should be used sparingly because they're not performant.

```
 1  from pyspark.sql.types import DateType,TimestampType
 2  from pyspark.sql.functions import udf
 3  from pyspark.sql.functions import to_date,to_timestamp
 4
 5  date_calc_udf = udf(funcDateCalc, DateType())
 6
 7  dfResult = dfTable1.withColumn("p_calc_ts, date_calc_udf(dfTable1['r
 8
 9  def funcDateCalc(inputDate, sendDate,  flag, addDays):
10
11      from datetime import datetime, timedelta
12
13      if(flag == 1):
14          if(sendDate is None):
15              modified_date = inputDate + timedelta(days=int(addDays))
16              return modified_date
17          else:
18              modified_date = sendDate + timedelta(days=int(addDays))
19              return modified_date
20      else:
21          return None
22
```

**Q10.** How do you write the Dataframe results to a relational database table?
**A10.**

```
 1
 2  dfResult.write \
 3      .mode('overwrite') \
 4      .format("jdbc") \
 5      .option("url", jdbcUrlWrite) \
 6      .option("dbtable", "my-table") \
 7      .option("user", "root") \
 8      .option("password", "password") \
 9      .option("batchsize", 10000) \
10      .option("driver", "com.mysql.jdbc.Driver") \
```

```
11        .option("truncate", 'false') \
12        .option("useSSL", false) \
13        .save()
14
```

**Q11.** How do you write the Dataframe results to a AWS s3 bucket?

**A11.**

```
1
2  dfResult.coalesce(1)
3        .write.csv("s3://my-bucket/some-path/output-file.csv", sep="|
4
```

**Q12.** What is a lit function?

**A12.** "lit" is a Spark SQL built-in function. It creates a Column of literal value.

```
1
2  from pyspark.sql.functions import lit
3  from datetime import date,datetime
4
5  send_date = datetime.now().date()
6  dfTable1.withColumn("run_date", lit(send_date ))
7
```

**Q13.** How do you create an SQLContext?

**A13.**

```
1
2  from pyspark.sql import SQLContext
3
4  spark = SparkSession.builder.appName("my-saprk")
5  sc = spark.sparkContext
6  sqlContext = SQLContext(sc)
7
```

**Q14.** How do you use sql like syntax in SparkSQL?

**A14.**

```
1
2  dfResult = spark.sql("SELECT * FROM dfTable1 WHERE NOT (calc_date  IS
3
```

**Note:** Refer to the pyspark SQL module API for more detail – http://spark.apache.org/docs/2.1.0/api/python/pyspark.sql.html.

Brush up with more examples? PySpark on Databricks examples.

Arulkumaran

## 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