Spark by {Examples}(https://sparkbyexamples.com/)m/)

Spark (https://sparkbyexamples.com/)

## Spark Tutorial

Spark – Installation on Windows (https://sparkbyexamples.com/spark/apache-spark-installation-on-windows/)

Spark – Installation on Linux | Ubuntu (https://sparkbyexamples.com/spark/spark-installation-on-linux-ubuntu/)

Spark – Cluster Setup with Hadoop Yarn (https://sparkbyexamples.com/spark/spark-setup-on-hadoop-yarn/)

Spark – Web/Application UI (https://sparkbyexamples.com/spark/spark-web-ui-understanding/)

Spark – Setup with Scala and IntelliJ (https://sparkbyexamples.com/spark/spark-setup-run-with-scala-intellij/)

Spark – How to Run Examples From this Site on IntelliJ IDEA (https://sparkbyexamples.com/spark/how-to-run-spark-examples-from-intellij/)

Spark – SparkSession (https://sparkbyexamples.com/spark/sparksession-explained-with-examples/)

Spark – SparkContext (https://sparkbyexamples.com/spark/spark-sparkcontext/)

## Spark RDD Tutorial

Spark RDD – Parallelize (https://sparkbyexamples.com/apache-spark-rdd/how-to-create-an-rdd-using-parallelize/)

PySpark   (https://sparkbyexamples.com/pyspark-tutorial/)

Hive   (https://sparkbyexamples.com/apache-hive-tutorial/)

**Prove your data skills**

Get the data science skills most companies are looking for. Invest in yourself today!

HBase   (https://sparkbyexamples.com/apache-hbase-tutorial/)

# Spark SQL – How to Remove Duplicate Rows

Kafka   (https://sparkbyexamples.com/apache-kafka-tutorials-with-examples/)

NNK (https://sparkbyexamples.com/author/admin/)  -  Apache Spark (https://sparkbyexamples.com/category/spark/)

FAQ's   (https://sparkbyexamples.com/spark-questions/)

More   (https://sparkbyexamples.com/)

Duplicate rows could be remove or drop from Spark SQL DataFrame using `distinct()` and `dropDuplicates()` functions, distinct() can be used to remove rows that have the same values on all columns whereas dropDuplicates() can be used to remove rows that have the same values on multiple selected columns.

Before we start, first let's create a
DataFrame
(https://sparkbyexamples.com/spark/diffe
rent-ways-to-create-a-spark-
dataframe/) with some duplicate rows
and duplicate values on a few columns.

```scala
import spark.implicits._

val simpleData = Seq(("James", "
  ("Michael", "Sales", 4600),
  ("Robert", "Sales", 4100),
  ("Maria", "Finance", 3000),
  ("James", "Sales", 3000),
  ("Scott", "Finance", 3300),
  ("Jen", "Finance", 3900),
  ("Jeff", "Marketing", 3000),
  ("Kumar", "Marketing", 2000),
  ("Saif", "Sales", 4100)
)
val df = simpleData.toDF("employ
df.show()
```

Yields below output

## Spark SQL Tutorial

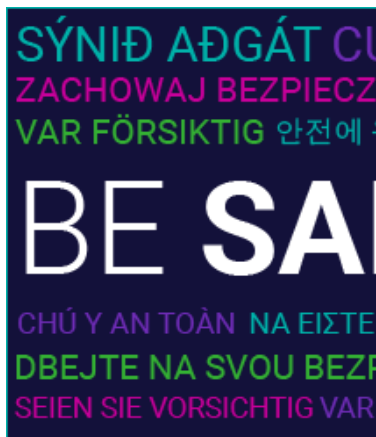```
+------------+----------+------
|employee_name|department|salary
+------------+----------+------
|      James|     Sales| 3000
|    Michael|     Sales| 4600
|     Robert|     Sales| 4100
|      Maria|   Finance| 3000
|      James|     Sales| 3000
|      Scott|   Finance| 3300
|        Jen|   Finance| 3900
|       Jeff| Marketing| 3000
|      Kumar| Marketing| 2000
|       Saif|     Sales| 4100
+------------+----------+------
```

On the above table, I've highlighted all duplicate rows, As you notice we have 2 rows that have duplicate values on all columns and we have 4 rows that have duplicate values on "department" and "salary" columns.

# 1. Use distinct() – Remove Duplicate Rows on DataFrame

On the above dataset, we have a total of 10 rows and one row with all values duplicated, performing distinct on this DataFrame should get us 9 as we have one duplicate.

```
//Distinct all columns
val distinctDF = df.distinct()
println("Distinct count: "+disti
distinctDF.show(false)
```

`distinct()` function on DataFrame returns a new DataFrame after removing the duplicate records. This example yields the below output.

```
Distinct count: 9
+-------------+----------+------
|employee_name|department|salary
+-------------+----------+------
|James        |Sales     |3000
|Michael      |Sales     |4600
|Maria        |Finance   |3000
|Robert       |Sales     |4100
|Saif         |Sales     |4100
|Scott        |Finance   |3300
|Jeff         |Marketing |3000
|Jen          |Finance   |3900
|Kumar        |Marketing |2000
+-------------+----------+------
```

Alternatively, you can also run `dropDuplicates()` function which return a new DataFrame with duplicate rows removed.

```
val df2 = df.dropDuplicates()
println("Distinct count: "+df2.c
df2.show(false)
```

## 2. Use dropDuplicate() – Remove Duplicate Rows on DataFrame

Spark doesn't have a distinct method that takes columns that should run distinct on however, Spark provides another signature of dropDuplicates() function which takes multiple columns to eliminate duplicates.

Note that calling dropDuplicates() on DataFrame returns a new DataFrame with duplicate rows removed.

```
//Distinct using dropDuplicates
val dropDisDF = df.dropDuplicate
println("Distinct count of depar
dropDisDF.show(false)
```

Yields below output. If you notice the output, It dropped 2 records that are duplicate.

```
Distinct count of department & s
+-------------+----------+------
|employee_name|department|salary
+-------------+----------+------
|Jen          |Finance   |3900
|Maria        |Finance   |3000
|Scott        |Finance   |3300
|Michael      |Sales     |4600
|Kumar        |Marketing |2000
|Robert       |Sales     |4100
|James        |Sales     |3000
|Jeff         |Marketing |3000
+-------------+----------+------
```

## 3. Source code – Remove Duplicate Rows

**Spark Streaming & Kafka**

```scala
package com.sparkbyexamples.spar

import org.apache.spark.sql.Spar
import org.apache.spark.sql.func

object SQLDistinct extends App {

  val spark: SparkSession = Spar
    .master("local[1]")
    .appName("SparkByExamples.co
    .getOrCreate()

  spark.sparkContext.setLogLevel

  import spark.implicits._

  val simpleData = Seq(("James",
    ("Michael", "Sales", 4600),
    ("Robert", "Sales", 4100),
    ("Maria", "Finance", 3000),
    ("James", "Sales", 3000),
    ("Scott", "Finance", 3300),
    ("Jen", "Finance", 3900),
    ("Jeff", "Marketing", 3000),
    ("Kumar", "Marketing", 2000)
    ("Saif", "Sales", 4100)
  )
  val df = simpleData.toDF("empl
  df.show()

  //Distinct all columns
  val distinctDF = df.distinct()
  println("Distinct count: "+dis
  distinctDF.show(false)

  val df2 = df.dropDuplicates()
  println("Distinct count: "+df2
  df2.show(false)

  //Distinct using dropDuplicate
  val dropDisDF = df.dropDuplica
  println("Distinct count of dep
  dropDisDF.show(false)

}
```

The complete example is available at GitHub (https://github.com/spark-examples/spark-scala-examples/blob/master/src/main/scala/com/sparkbyexamples/spark/dataframe/functions/aggregate/DistinctCount.scala) for reference.

## 4. Conclusion

In this Spark article, you have learned how to remove DataFrame rows that are exact duplicates using `distinct()` and learned how to remove duplicate rows based on multiple columns using `dropDuplicate()` function with Scala example.

**Share this:**

**TAGS: DISTINCT() (HTTPS://SPARKBYEXAMPLES.COM/TAG/DISTINCT/) , DROPDUPLICATES() (HTTPS://SPARKBYEXAMPLES.COM/TAG/DROPDUPLICATES/)**

**NNK (Https://Sparkbyexamples.Com/Author/Admin/)**

(https://sp

**Implementation Specialists**

Hiring for multiple locations in India

Hexaware Technolog

SparkByExamples.com is a Big Data and Spark examples community page, all examples are simple and easy to understand and well tested in our development environment Read more .. (https://sparkbyexamples.com/about-sparkbyexamples/)

> **THIS POST HAS 4 COMMENTS**

**Anonymous**

10 APR 2021     REPLY

Awesome

**Anonymous**

13 MAR 2021     REPLY
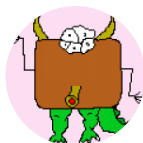
These tutorials are realy helpful, it would have been very good, if written in pyhton

**NNK**                     13 MAR 2021

Hi, There are many articles wri Python as-well. Please refer to

Remove duplicate rows in PySp (Spark with Python) (https://sparkbyexamples.com/p k/pyspark-distinct-to-drop-duplicates/)

**Anonymous**

8 MAR 2021     REPLY

Thank you, that helped.

## Leave a Reply

## Overview of digital pathology

### AI use in image analysis

Find out more about AI, how to use it in digital pathology, with our free resource

## Categories

Apache Hadoop (https://sparkbyexamples.com/category/hadoop/)

Apache Spark (https://sparkbyexamples.com/category/spark/)

Apache Spark Streaming (https://sparkbyexamples.com/category/spark/apache-spark-streaming/)

Apache Kafka (https://sparkbyexamples.com/category/kafka/)

Apache HBase (https://sparkbyexamples.com/categor

## Recent Posts

Spark regexp_replace() – Replace String Value (https://sparkbyexamples.com/spark-regexp_replace-replace-string-value/)

How to Run a PySpark Script from Python? (https://sparkbyexamples.com/pyspark/run-pyspark-script-from-python-subprocess/)

Spark SQL like() Using Wildcard Example (https://sparkbyexamples.com/spark/spark-sql-like-using-wildcard-example/)

Spark isin() & IS NOT IN Operator Example

## About SparkByExamples.Com

SparkByExamples.com is a Big Data and Spark examples community page, all examples are simple and easy to understand, and well tested in our development environment Read more .. (https://sparkbyexamples.com/about-sparkbyexamples/)

## Follow Us

(https:  (https:

//www.  //www.

(https:  facebo  linkedi  (https:

//twitte  ok.co  n.com/  //githu

r.com/  m/spar  in/n-  b.com/

sparkb  kbyex  nk-  spark-

yexam  ample  b860a  examp

ples)  s/)  8193/)  les/)