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

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

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

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

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

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

More ∨   (https://sparkbyexamples.com/)

# PySpark – Distinct to Drop Duplicate Rows

👤 NNK (https://sparkbyexamples.com/author/admin/) -
🗀 PySpark (https://sparkbyexamples.com/category/pyspark/)

PySpark `distinct()` function is used to drop/remove the duplicate rows (all columns) from DataFrame and `dropDuplicates()` is used to drop rows based on selected (one or multiple) columns. In this article, you will learn how to use distinct() and dropDuplicates() functions with PySpark example.

Before we start, first let's create a DataFrame (https://sparkbyexamples.com/pyspark/different-ways-to-create-dataframe-in-pyspark/) with some duplicate rows and

values on a few columns. We use this
DataFrame to demonstrate how to get
distinct multiple columns.

```python
import pyspark
from pyspark.sql import SparkSe
from pyspark.sql.functions impo
spark = SparkSession.builder.ap

data = [("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) \
  ]
columns= ["employee_name", "dep
df = spark.createDataFrame(data
df.printSchema()
df.show(truncate=False)
```

Yields below output

```
+------------+---------+-----
|employee_name|department|salar
+------------+---------+-----
|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, record with
employer name Robert has 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. Get Distinct Rows (By Comparing All Columns)

On the above DataFrame, we have a total of 10 rows with 2 rows having all values duplicated, performing distinct on this DataFrame should get us 9 after removing 1 duplicate row.

```
distinctDF = df.distinct()
print("Distinct count: "+str(di
distinctDF.show(truncate=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|salar
+------------+----------+-----
|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 returns a new DataFrame after removing duplicate rows.

```
df2 = df.dropDuplicates()
print("Distinct count: "+str(df
df2.show(truncate=False)
```

## 2. PySpark Distinct of Selected Multiple Columns

PySpark doesn't have a distinct method which takes columns that should run distinct on (drop duplicate rows on selected multiple columns) however, it 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.

```
dropDisDF = df.dropDuplicates([
print("Distinct count of depart
dropDisDF.show(truncate=False)
```
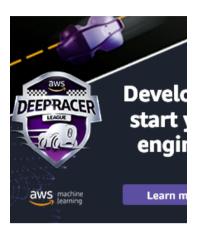
Yields below output. If you notice the output, It dropped 2 records that are duplicate.

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

## Source Code to Get

## Distinct Rows

## PySpark Built-In Functions

```python
import pyspark
from pyspark.sql import SparkSe
from pyspark.sql.functions impo
spark = SparkSession.builder.ap

data = [("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) \
  ]
columns= ["employee_name", "dep
df = spark.createDataFrame(data
df.printSchema()
df.show(truncate=False)

#Distinct
distinctDF = df.distinct()
print("Distinct count: "+str(di
distinctDF.show(truncate=False)

#Drop duplicates
df2 = df.dropDuplicates()
print("Distinct count: "+str(df
df2.show(truncate=False)

#Drop duplicates on selected co
dropDisDF = df.dropDuplicates([
print("Distinct count of depart
dropDisDF.show(truncate=False)
}
```

The complete example is available at
GitHub (https://github.com/spark-
examples/pyspark-
examples/blob/master/pyspark-
distinct.py) for reference.

## Conclusion

In this PySpark SQL article, you have
learned distinct() method which is
used to get the distinct values of rows
(all columns) and also learned how to
use dropDuplicates() to get the
distinct and finally learned using
dropDuplicates() function to get distinct
of multiple columns.

Happy Learning !!

## Related Articles:

- 

---

**Share this:**

**TAGS:** **DISTINCT ROWS
(HTTPS://SPARKBYEXAMPLES.COM/TAG/DISTINCT
-ROWS/)**, **DISTINCT()
(HTTPS://SPARKBYEXAMPLES.COM/TAG/DISTINCT
/)**, **DROPDUPLICATES()
(HTTPS://SPARKBYEXAMPLES.COM/TAG/DROPDU
PLICATES/)**, **DUPLICATE ROWS
(HTTPS://SPARKBYEXAMPLES.COM/TAG/DUPLICA
TE-ROWS/)**

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

(https://sp
arkbyexa
mples.co
m/author/
admin/)

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**

**abdulsattar**

4 APR 2021          REPLY

bro please correct them above if there is no duplicate , i spent alot time on this and after all i came to see if there is someone else having same problem

**NNK**   4 APR 2021          REPLY

Hi Abdulsattar, I have updated the article when it was pointed out the first time. You should not see duplicate() function used anywhere. Could you please let me know where you are seeing, may be it's caching somewhere.

**Sneha**   31 OCT 2020          REPLY

I don't see duplicate() method used, is there a confusion between distinct() and duplicate() ? Please check.

**NNK**   1 NOV 2020          REPLY

Thanks Sneha. Yes, it should be distinct(), there is no duplicate but PySpark also has dropDuplicates().

## Leave a Reply

## Categories

Apache Hadoop
(https://sparkbyexamples.com/catego
ry/hadoop/)

Apache Spark
(https://sparkbyexamples.com/catego
ry/spark/)

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

Apache Kafka
(https://sparkbyexamples.com/catego
ry/kafka/)

Apache HBase
(https://sparkbyexamples.com/catego
ry/hbase/)

Apache Cassandra
(https://sparkbyexamples.com/catego
ry/cassandra/)

Snowflake Database
(https://sparkbyexamples.com/catego

## Recent Posts

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

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

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

Spark isin() & IS NOT IN Operator
Example
(https://sparkbyexamples.com/spark/sp
ark-isin-is-not-in-operator-example/)

Spark – Get Size/Length of Array & Map
Column
(https://sparkbyexamples.com/spark/sp
ark-get-size-length-of-array-map-
column/)

## 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

Spark Using Length/Size Of a
DataFrame Column
(https://sparkbyexamples.com/spark/sp
ark-using-length-size-of-a-dataframe-
column/)

Spark rlike() Working with Regex
Matching Examples
(https://sparkbyexamples.com/spark/sp
ark-rlike-regex-matching-examples/)

Spark Check String Column Has
Numeric Values
(https://sparkbyexamples.com/spark/sp
ark-check-string-column-has-numeric-
values/)

Spark Check Column Data Type is
Integer or String
(https://sparkbyexamples.com/spark/sp
ark-check-column-data-type-is-integer-
or-string/)

(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/)