Spark by {Examples} (https://sparkbyexamples.com/) Spark + (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 Where Filter Function | Multiple Conditions

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

PySpark `filter()` function is used to
filter the rows from RDD/DataFrame
based on the given condition or SQL
expression, you can also use `where()`
clause instead of the filter() if you are
coming from an SQL background, both
these functions operate exactly the
same.

In this PySpark article, you will learn
how to apply a filter on DataFrame
columns of string, arrays, struct types
by using single and multiple conditions
and also applying filter using `isin()`
with PySpark (Python Spark) examples.

**Related Article:**

- How to Filter Rows with NULL/NONE
  (IS NULL & IS NOT NULL) in
  PySpark
  (https://sparkbyexamples.com/pyspar
  k/pyspark-filter-rows-with-null-
  values/)

**Note:** PySpark Column Functions (https://sparkbyexamples.com/pyspark/pyspark-column-functions/) provides several options that can be used with filter().

# 1. PySpark DataFrame filter() Syntax

Below is syntax of the filter function. condition would be an expression you wanted to filter.

```
filter(condition)
```

Before we start with examples, first let's create a DataFrame (https://sparkbyexamples.com/pyspark/different-ways-to-create-dataframe-in-pyspark/). Here, I am using a DataFrame with StructType (https://sparkbyexamples.com/pyspark/pyspark-structtype-and-structfield/) and ArrayType (https://sparkbyexamples.com/pyspark/pyspark-arraytype-column-with-examples/) columns as I will also be covering examples with struct and array types as-well.

```python
from pyspark.sql.types import S
from pyspark.sql.types import S
data = [
    (("James","","Smith"),["Jav
    (("Anna","Rose",""),["Spark
    (("Julia","","Williams"),["
    (("Maria","Anne","Jones"),[
    (("Jen","Mary","Brown"),["C
    (("Mike","Mary","Williams")
 ]

schema = StructType([
    StructField('name', Struct
        StructField('firstname'
        StructField('middlename
         StructField('lastname'
    ])),
    StructField('languages', A
    StructField('state', Strin
    StructField('gender', Stri
 ])

df = spark.createDataFrame(data
df.printSchema()
df.show(truncate=False)
```

This yields below schema and
DataFrame results.

```
root
 |-- name: struct (nullable = t
 |    |-- firstname: string (nu
 |    |-- middlename: string (n
 |    |-- lastname: string (nul
 |-- languages: array (nullable
 |    |-- element: string (cont
 |-- state: string (nullable =
 |-- gender: string (nullable =

+--------------------+-------
|name                |languag
+--------------------+-------
|[James, , Smith]    |[Java,
|[Anna, Rose, ]      |[Spark,
|[Julia, , Williams] |[CSharp
|[Maria, Anne, Jones]|[CSharp
|[Jen, Mary, Brown]  |[CSharp
|[Mike, Mary, Williams]|[Python
+--------------------+-------
```

# 2. DataFrame filter() with Column Condition

Use Column with the condition to filter the rows from DataFrame, using this you can express complex condition by referring column names using `dfObject.colname`

```
# Using equals condition
df.filter(df.state == "OH").sho

+--------------------+-------
|name                |languag
+--------------------+-------
|[James, , Smith]    |[Java,
|[Julia, , Williams] |[CSharp
|[Mike, Mary, Williams]|[Python
+--------------------+-------

# not equals condition
df.filter(df.state != "OH") \
    .show(truncate=False)
df.filter(~(df.state == "OH"))
    .show(truncate=False)
```

Same example can also written as below. In order to use this first you need to import `from pyspark.sql.functions import col`
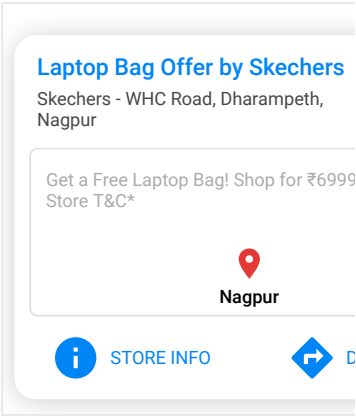
```
#Using SQL col() function
from pyspark.sql.functions impo
df.filter(col("state") == "OH")
    .show(truncate=False)
```

# 3. DataFrame filter() with SQL Expression

If you are coming from SQL background, you can use that knowledge in PySpark to filter DataFrame rows with SQL expressions.

```
Jsing SQL Expression
F.filter("gender == 'M'").show
For not equal
F.filter("gender != 'M'").show
F.filter("gender <> 'M'").show
```

# 4. PySpark Filter with Multiple Conditions

In PySpark, to `filter()` rows on DataFrame based on multiple conditions, you case use either `Column` with a condition or SQL expression. Below is just a simple example using AND (&) condition, you can extend this with OR(|), and NOT(!) conditional expressions as needed.

```
//Filter multiple condition
df.filter( (df.state  == "OH")
    .show(truncate=False)
```

This yields below DataFrame results.

```
+--------------------+-------
|name                |languag
+--------------------+-------
|[James, , Smith]    |[Java,
|[Mike, Mary, Williams]|[Python
+--------------------+-------
```

# 5. Filter Based on List Values

If you have a list of elements and you wanted to filter that is not in the list or in the list, use `isin()` function of Column class (https://sparkbyexamples.com/pyspark/pyspark-column-functions/) and it doesn't have isnotin() function but you do the same using not operator (~)

```
#Filter IS IN List values
li=["OH","CA","DE"]
df.filter(df.state.isin(li)).sh
+-------------------+---------
|               name|
+-------------------+---------
|    [James, , Smith]|[Java, Sc
| [Julia, , Williams]|      [CS
|[Mike, Mary, Will...|      [Py
+-------------------+---------

# Filter NOT IS IN List values
#These show all records with NY
df.filter(~df.state.isin(li)).s
df.filter(df.state.isin(li)==Fa
```

## 6. Filter Based on Starts With, Ends With, Contains

You can also filter DataFrame rows by using `startswith()`, `endswith()` and `contains()` methods of Column class. For more examples on Column class, refer to PySpark Column Functions (https://sparkbyexamples.com/pyspark/pyspark-column-functions/).

```
# Using startswith
df.filter(df.state.startswith("I
+-------------------+---------
|               name|
+-------------------+---------
|      [Anna, Rose, ]|[Spark, J
|[Maria, Anne, Jones]|      [CS
|   [Jen, Mary, Brown]|      [CS
+-------------------+---------

#using endswith
df.filter(df.state.endswith("H"

#contains
df.filter(df.state.contains("H"
```

## 7. PySpark Filter like and rlike

If you have SQL background you must be familiar with `like` and `rlike` (regex like), PySpark also provides similar methods in Column class to filter similar values using wildcard characters. You can use rlike() to filter by checking values case insensitive.

```
data2 = [(2,"Michael Rose"),(3,
     (4,"Rames Rose"),(5,"Rames
  ]
df2 = spark.createDataFrame(dat

# like - SQL LIKE pattern
df2.filter(df2.name.like("%rose
+---+----------+
| id|      name|
+---+----------+
|  5|Rames rose|
+---+----------+

# rlike - SQL RLIKE pattern (LI
#This check case insensitive
df2.filter(df2.name.rlike("(?i)
+---+------------+
| id|        name|
+---+------------+
|  2|Michael Rose|
|  4|  Rames Rose|
|  5|  Rames rose|
```

## 8. Filter on an Array column

When you want to filter rows from DataFrame based on value present in an array collection column, you can use the first syntax. The below example uses `array_contains()` from Pyspark SQL functions (https://sparkbyexamples.com/spark/spark-sql-functions/) which checks if a value contains in an array if present it returns true otherwise false.

```python
from pyspark.sql.functions impo
df.filter(array_contains(df.lan
    .show(truncate=False)
```

This yields below DataFrame results.

```
+----------------+------------
|name            |languages
+----------------+------------
|[James, , Smith]|[Java, Scala,
|[Anna, Rose, ]  |[Spark, Java,
+----------------+------------
```

## 9. Filtering on Nested Struct columns

If your DataFrame consists of nested struct columns, you can use any of the above syntaxes to filter the rows based on the nested column.

```python
//Struct condition
df.filter(df.name.lastname == "
    .show(truncate=False)
```

This yields below DataFrame results

```
+--------------------+-------
|name                |languag
+--------------------+-------
|[Julia, , Williams] |[CSharp
|[Mike, Mary, Williams]|[Python
+--------------------+-------
```

## 10. Source code of PySpark where filter

```python
import pyspark
from pyspark.sql import SparkSe
from pyspark.sql.types import S
from pyspark.sql.functions impo

spark = SparkSession.builder.ap

arrayStructureData = [
        (("James","","Smith"),[
        (("Anna","Rose",""),["S
        (("Julia","","Williams"
        (("Maria","Anne","Jones
        (("Jen","Mary","Brown")
        (("Mike","Mary","Willia
        ]

arrayStructureSchema = StructTy
        StructField('name', Str
            StructField('first
            StructField('middl
            StructField('lastn
            ])),
        StructField('languages
        StructField('state', S
        StructField('gender',
        ])


df = spark.createDataFrame(data
df.printSchema()
df.show(truncate=False)

df.filter(df.state == "OH") \
    .show(truncate=False)

df.filter(col("state") == "OH")
    .show(truncate=False)

df.filter("gender  == 'M'") \
    .show(truncate=False)

df.filter( (df.state  == "OH")
    .show(truncate=False)

df.filter(array_contains(df.lan
    .show(truncate=False)

df.filter(df.name.lastname == "
    .show(truncate=False)
```

Examples explained here are also
available at PySpark examples GitHub
(https://github.com/spark-

examples/pyspark-
examples/blob/master/pyspark-filter.py)
project for reference.

## 11. Conclusion

In this tutorial, I've explained how to
filter rows from PySpark DataFrame
based on single or multiple conditions
and SQL expression, also learned
filtering rows by providing conditions on
the array and struct column with Spark
with Python examples.

Alternatively, you can also use
`where()` function to filter the rows on
PySpark DataFrame.

Happy Learning !!

---

**Share this:**

**TAGS:** **FILTER()
(HTTPS://SPARKBYEXAMPLES.COM/TAG/FILTER/)**,
**WHERE()
(HTTPS://SPARKBYEXAMPLES.COM/TAG/WHERE/)**

---

**NNK**

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

(https://sp
arkbyexa
mples.co

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 ONE COMMENT**

**Anonymous**

**24 MAR 2021**      **REPLY**

I am new to pyspark and
this blog was extremely
helpful to understand the
concept. Thank you!!

## Leave a Reply

Categories

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

Recent Posts

Spark regexp_replace() – Replace
String Value
(https://sparkbyexamples.com/spark/sp

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/)