

# PySpark Filter() Function

PySpark `filter()` function is used to filter the rows from RDD/DataFrame based on the given condition or SQL expression, we can also use `where()` clause instead of the `filter()`. Both these functions operate the same. `filter()` function returns a new DataFrame or RDD with only the rows that meet the condition specified.

Before we start with examples, first let's [create a DataFrame](#).

```
1  from pyspark.sql.types import StructType, StructField
2  from pyspark.sql.types import StringType, IntegerType, ArrayType
3  data = [
4      ("James", "", "Smith"), ["Java", "Scala", "C++"], "OH", "M"),
5      ("Anna", "Rose", ""), ["Spark", "Java", "C++"], "NY", "F"),
6      ("Julia", "", "Williams"), ["CSharp", "VB"], "OH", "F"),
7      ("Maria", "Anne", "Jones"), ["CSharp", "VB"], "NY", "M"),
8      ("Jen", "Mary", "Brown"), ["CSharp", "VB"], "NY", "M"),
9      ("Mike", "Mary", "Williams"), ["Python", "VB"], "OH", "M")
10 ]
11
12 schema = StructType([
13     StructField('name', StructType([
14         StructField('firstname', StringType(), True),
15         StructField('middlename', StringType(), True),
16         StructField('lastname', StringType(), True)
17     ])),
18     StructField('languages', ArrayType(StringType()), True),
19     StructField('state', StringType(), True),
20     StructField('gender', StringType(), True)
21 ])
22
23 df = spark.createDataFrame(data = data, schema = schema)
24 df.printSchema()
25 df.show(truncate=False)
```

```

root
|-- name: struct (nullable = true)
|   |-- firstname: string (nullable = true)
|   |-- middlename: string (nullable = true)
|   |-- lastname: string (nullable = true)
|-- languages: array (nullable = true)
|   |-- element: string (containsNull = true)
|-- state: string (nullable = true)
|-- gender: string (nullable = true)

+-----+-----+-----+-----+
|name           |languages           |state|gender|
+-----+-----+-----+-----+
|{James, , Smith}|[Java, Scala, C++]|OH   |M     |
|{Anna, Rose, }  |[Spark, Java, C++]|NY   |F     |
|{Julia, , Williams}|[CSharp, VB]      |OH   |F     |
|{Maria, Anne, Jones}|[CSharp, VB]      |NY   |M     |
|{Jen, Mary, Brown}|[CSharp, VB]      |NY   |M     |
|{Mike, Mary, Williams}|[Python, VB]      |OH   |M     |
+-----+-----+-----+-----+

```

## DataFrame filter() with Column Condition

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

► (1) Spark Jobs

```

+-----+-----+-----+-----+
|name           |languages           |state|gender|
+-----+-----+-----+-----+
|{James, , Smith}|[Java, Scala, C++]|OH   |M     |
|{Julia, , Williams}|[CSharp, VB]      |OH   |F     |
|{Mike, Mary, Williams}|[Python, VB]      |OH   |M     |
+-----+-----+-----+-----+

```

```

1 df.filter(df.state != "OH") \
2 | .show(truncate=False)

```

► (1) Spark Jobs

```

+-----+-----+-----+-----+
|name          |languages          |state|gender|
+-----+-----+-----+-----+
|{Anna, Rose, }|[Spark, Java, C++]|NY   |F     |
|{Maria, Anne, Jones}|[CSharp, VB]      |NY   |M     |
|{Jen, Mary, Brown}|[CSharp, VB]      |NY   |M     |
+-----+-----+-----+-----+

```

```

1 df.filter(~(df.state == "OH")) \
2 | .show(truncate=False)

```

► (1) Spark Jobs

```

+-----+-----+-----+-----+
|name          |languages          |state|gender|
+-----+-----+-----+-----+
|{Anna, Rose, }|[Spark, Java, C++]|NY   |F     |
|{Maria, Anne, Jones}|[CSharp, VB]      |NY   |M     |
|{Jen, Mary, Brown}|[CSharp, VB]      |NY   |M     |
+-----+-----+-----+-----+

```

```

1  from pyspark.sql.functions import col
2  df.filter(col("state") == "OH") \
3  |    .show(truncate=False)

```

► (1) Spark Jobs

name	languages	state	gender
{James, , Smith}	[Java, Scala, C++]	OH	M
{Julia, , Williams}	[CSharp, VB]	OH	F
{Mike, Mary, Williams}	[Python, VB]	OH	M

## DataFrame filter() with SQL Expression

```

1  df.filter("gender == 'M'").show()

```

► (1) Spark Jobs

name	languages	state	gender
{James, , Smith}	[Java, Scala, C++]	OH	M
{Maria, Anne, Jones}	[CSharp, VB]	NY	M
{Jen, Mary, Brown}	[CSharp, VB]	NY	M
{Mike, Mary, Will...}	[Python, VB]	OH	M

```
1 df.filter("gender != 'M'").show()
```

► (1) Spark Jobs

name	languages	state	gender
{Anna, Rose, }	[Spark, Java, C++]	NY	F
{Julia, , Williams}	[CSharp, VB]	OH	F

```
1 df.filter("gender <> 'M'").show()
```

► (1) Spark Jobs

name	languages	state	gender
{Anna, Rose, }	[Spark, Java, C++]	NY	F
{Julia, , Williams}	[CSharp, VB]	OH	F

## PySpark Filter with Multiple Conditions

```
1 df.filter( (df.state == "OH") & (df.gender == "M") ) \  
2 | .show(truncate=False)
```

► (1) Spark Jobs

name	languages	state	gender
{James, , Smith}	[Java, Scala, C++]	OH	M
{Mike, Mary, Williams}	[Python, VB]	OH	M

## Filter Based on List Values

```
1 li=["OH","CA","DE"]  
2 df.filter(df.state.isin(li)).show()
```

► (1) Spark Jobs

name	languages	state	gender
{James, , Smith}	[Java, Scala, C++]	OH	M
{Julia, , Williams}	[CSharp, VB]	OH	F
{Mike, Mary, Will...}	[Python, VB]	OH	M

```

1  li=["OH","CA","DE"]
2  df.filter(~df.state.isin(li)).show()

```

► (1) Spark Jobs

name	languages	state	gender
{Anna, Rose, }	[Spark, Java, C++]	NY	F
{Maria, Anne, Jones}	[CSharp, VB]	NY	M
{Jen, Mary, Brown}	[CSharp, VB]	NY	M

```

1  li=["OH","CA","DE"]
2  df.filter(df.state.isin(li)==False).show()

```

► (1) Spark Jobs

name	languages	state	gender
{Anna, Rose, }	[Spark, Java, C++]	NY	F
{Maria, Anne, Jones}	[CSharp, VB]	NY	M
{Jen, Mary, Brown}	[CSharp, VB]	NY	M

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

```
1 # Using startswith
2 df.filter(df.state.startswith("N")).show()
```

► (1) Spark Jobs

name	languages	state	gender
{Anna, Rose, }	[Spark, Java, C++]	NY	F
{Maria, Anne, Jones}	[CSharp, VB]	NY	M
{Jen, Mary, Brown}	[CSharp, VB]	NY	M

```
1 #using endswith
2 df.filter(df.state.endswith("H")).show()
```

► (1) Spark Jobs

name	languages	state	gender
{James, , Smith}	[Java, Scala, C++]	OH	M
{Julia, , Williams}	[CSharp, VB]	OH	F
{Mike, Mary, Will...}	[Python, VB]	OH	M



```

1  #contains
2  df.filter(df.state.contains("H")).show()

```

► (1) Spark Jobs

name	languages	state	gender
{James, , Smith}	[Java, Scala, C++]	OH	M
{Julia, , Williams}	[CSharp, VB]	OH	F
{Mike, Mary, Will...}	[Python, VB]	OH	M

## PySpark Filter “like”

```

1  # like - SQL LIKE pattern
2  df.filter(df.name.firstname.like("%J%")).show()

```

► (1) Spark Jobs

name	languages	state	gender
{James, , Smith}	[Java, Scala, C++]	OH	M
{Julia, , Williams}	[CSharp, VB]	OH	F
{Jen, Mary, Brown}	[CSharp, VB]	NY	M

## Filter on an Array column

```
1 from pyspark.sql.functions import array_contains
2 df.filter(array_contains(df.languages, "Java")) \
3 | .show(truncate=False)
```

### ► (2) Spark Jobs

```
+-----+-----+-----+-----+
|name          |languages          |state|gender|
+-----+-----+-----+-----+
|{James, , Smith}|[Java, Scala, C++]|OH   |M     |
|{Anna, Rose, }  |[Spark, Java, C++]|NY   |F     |
+-----+-----+-----+-----+
```

## Filtering on Nested Struct columns

```
1 # Struct condition
2 df.filter(df.name.lastname == "Williams") \
3 | .show(truncate=False)
```

### ► (1) Spark Jobs

```
+-----+-----+-----+-----+
|name          |languages          |state|gender|
+-----+-----+-----+-----+
|{Julia, , Williams}|[CSharp, VB]|OH   |F     |
|{Mike, Mary, Williams}|[Python, VB]|OH   |M     |
+-----+-----+-----+-----+
```