# Spark by {Examples} (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/)

Spark RDD – Read text file (https://sparkbyexamples.com/apache-spark-rdd/spark-read-multiple-text-files-into-a-single-rdd/)

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/)

-30%    -30%    -40%

shop.adidas.co.in

# Spark Groupby Example with DataFrame

NNK (https://sparkbyexamples.com/author/admin/) -

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

Similar to SQL "GROUP BY" clause, Spark groupBy() function is used to collect the identical data into groups on DataFrame/Dataset and perform aggregate functions on the grouped data. In this article, I will explain several groupBy() examples with the Scala language.

The same approach can be used with the Pyspark (Spark with Python).

**Syntax:**

```
groupBy(col1 : scala.Predef.Str
        org.apache.spark.sql.Rela
```

When we perform `groupBy()` on Spark
Dataframe, it returns
`RelationalGroupedDataset` object
which contains below aggregate
functions.

`count()` - Returns the count of rows
for each group.

`mean()` - Returns the mean of values
for each group.

`max()` - Returns the maximum of
values for each group.

`min()` - Returns the minimum of
values for each group.

`sum()` - Returns the total for values
for each group.

`avg()` - Returns the average for
values for each group.

`agg()` - Using agg() function, we can
calculate more than one aggregate at a
time.

`pivot()` - This function is used to Pivot the DataFrame which I will not be covered in this article as I already have a dedicated article for Pivot & Unvot DataFrame (https://sparkbyexamples.com/spark/how-to-pivot-table-and-unpivot-a-spark-dataframe/).

# Preparing Data & DataFrame

Before we start, let's create the DataFrame (https://sparkbyexamples.com/spark/different-ways-to-create-a-spark-dataframe/) from a sequence of the data to work with. This DataFrame contains columns "`employee_name`", "`department`", "`state`", "`salary`", "`age`" and "`bonus`" columns.

We will use this Spark DataFrame to run groupBy() on "department" columns and calculate aggregates like minimum, maximum, average, total salary for each group using min(), max() and sum() aggregate functions respectively. and finally, we will also see how to do group and aggregate on multiple columns.

```
import spark.implicits._
val simpleData = Seq(("James","
    ("Michael","Sales","NY",860
    ("Robert","Sales","CA",8100
    ("Maria","Finance","CA",900
    ("Raman","Finance","CA",990
    ("Scott","Finance","NY",830
    ("Jen","Finance","NY",79000
    ("Jeff","Marketing","CA",80
    ("Kumar","Marketing","NY",9
  )
val df = simpleData.toDF("emplo
df.show()
```

Yields below output.

```
+------------+----------+-----
|employee_name|department|state
+------------+----------+-----
|      James|     Sales|    NY
|    Michael|     Sales|    NY
|     Robert|     Sales|    CA
|      Maria|   Finance|    CA
|      Raman|   Finance|    CA
|      Scott|   Finance|    NY
|        Jen|   Finance|    NY
|       Jeff| Marketing|    CA
|      Kumar| Marketing|    NY
+------------+----------+-----
```

# groupBy and aggregate on DataFrame columns

Let's do the `groupBy()` on `department` column of DataFrame and then find the sum of salary for each department using `sum()` aggregate function.

```
df.groupBy("department").sum("s
+----------+----------+
|department|sum(salary)|
+----------+----------+
|Sales     |257000    |
|Finance   |351000    |
|Marketing |171000    |
+----------+----------+
```

Similarly, we can calculate the number of employee in each department using `count()`

```
df.groupBy("department").count(
```

Calculate the minimum salary of each department using `min()`

```
df.groupBy("department").min("s
```

Calculate the maximin salary of each department using `max()`

```
df.groupBy("department").max("s
```

Calculate the average salary of each department using `avg()`

```
df.groupBy("department").avg( "
```

Calculate the mean salary of each department using `mean()`

```
df.groupBy("department").mean(
```

# groupBy and aggregate on multiple DataFrame columns

Similarly, we can also run groupBy and aggregate on two or more DataFrame columns, below example does group by on `department`, `state` and does sum() on `salary` and `bonus` columns.

```
//GroupBy on multiple columns
df.groupBy("department","state"
    .sum("salary","bonus")
    .show(false)
```

This yields the below output.

```
+----------+-----+-----------+-
|department|state|sum(salary)|s
+----------+-----+-----------+-
|Finance   |NY   |162000     |3
|Marketing |NY   |91000      |2
|Sales     |CA   |81000      |2
|Marketing |CA   |80000      |1
|Finance   |CA   |189000     |4
|Sales     |NY   |176000     |3
+----------+-----+-----------+-
```

similarly, we can run group by and
aggregate on tow or more columns for
other aggregate functions, please refer
below source code for example.

# Running more aggregates at a time

Using `agg()` aggregate function we
can calculate many aggregations at a
time on a single statement using Spark
SQL aggregate
(https://sparkbyexamples.com/spark/sp
ark-sql-aggregate-functions/) functions
sum()
(https://sparkbyexamples.com/spark/sp
ark-sql-aggregate-functions/#sum),
avg()
(https://sparkbyexamples.com/spark/sp
ark-sql-aggregate-functions/#avg),
min()
(https://sparkbyexamples.com/spark/sp
ark-sql-aggregate-functions/#min),
max()
(https://sparkbyexamples.com/spark/sp
ark-sql-aggregate-functions/#max)
mean()
(https://sparkbyexamples.com/spark/sp
ark-sql-aggregate-functions/#mean)
e.t.c. In order to use these, we should
import "`import
org.apache.spark.sql.functions.
_`"

```
import org.apache.spark.sql.fun
df.groupBy("department")
    .agg(
      sum("salary").as("sum_sal
      avg("salary").as("avg_sal
      sum("bonus").as("sum_bonu
      max("bonus").as("max_bonu
    .show(false)
```

This example does group on
`department` column and calculates
`sum()` and `avg()` of `salary` for each
department and calculates `sum()` and
`max()` of bonus for each department.

```
+----------+----------+--------
|department|sum_salary|avg_sala
+----------+----------+--------
|Sales     |257000    |85666.66
|Finance   |351000    |87750.0
|Marketing |171000    |85500.0
+----------+----------+--------
```

# Using filter on aggregate data

Similar to SQL "HAVING" clause, On
Spark DataFrame we can use either
where()
(https://sparkbyexamples.com/spark/wo
rking-with-spark-dataframe-filter/) or
filter()
(https://sparkbyexamples.com/spark/wo
rking-with-spark-dataframe-filter/)
function to filter the rows of aggregated
a.

```
f.groupBy("department")
    .agg(
      sum("salary").as("sum_sal
      avg("salary").as("avg_sal
      sum("bonus").as("sum_bonu
      max("bonus").as("max_bonu
    .where(col("sum_bonus") >=
    .show(false)
```

removes the sum of a bonus that
less than 50000 and yields below
out.

```
-----------+----------+--------
department|sum_salary|avg_sala
-----------+----------+--------
Sales      |257000    |85666.66
Finance    |351000    |87750.0
-----------+----------+--------
```

▶

urce code

```scala
package com.sparkbyexamples.spa

import org.apache.spark.sql.Spa
import org.apache.spark.sql.fun

object GroupbyExample extends A
  val spark: SparkSession = Spa
    .master("local[1]")
    .appName("SparkByExamples.c
    .getOrCreate()

  spark.sparkContext.setLogLeve
  import spark.implicits._

  val simpleData = Seq(("James"
    ("Michael","Sales","NY",860
    ("Robert","Sales","CA",8100
    ("Maria","Finance","CA",900
    ("Raman","Finance","CA",990
    ("Scott","Finance","NY",830
    ("Jen","Finance","NY",79000
    ("Jeff","Marketing","CA",80
    ("Kumar","Marketing","NY",9
  )
  val df = simpleData.toDF("emp
  df.show()

  //Group By on single column
  df.groupBy("department").coun
  df.groupBy("department").avg(
  df.groupBy("department").sum(
  df.groupBy("department").min(
  df.groupBy("department").max(
  df.groupBy("department").mean

  //GroupBy on multiple columns
  df.groupBy("department","stat
    .sum("salary","bonus")
    .show(false)
  df.groupBy("department","stat
    .avg("salary","bonus")
    .show(false)
  df.groupBy("department","stat
    .max("salary","bonus")
    .show(false)
  df.groupBy("department","stat
    .min("salary","bonus")
    .show(false)
  df.groupBy("department","stat
    .mean("salary","bonus")
    .show(false)

  //Running Filter
  df.groupBy("department","stat
    .sum("salary","bonus")
```

```
        .show(false)

  //using agg function
  df.groupBy("department")
    .agg(
      sum("salary").as("sum_sal
      avg("salary").as("avg_sal
      sum("bonus").as("sum_bonu
      max("bonus").as("max_bonu
    .show(false)

  df.groupBy("department")
    .agg(
      sum("salary").as("sum_sal
      avg("salary").as("avg_sal
      sum("bonus").as("sum_bonu
      stddev("bonus").as("stdde
    .where(col("sum_bonus") > 5
    .show(false)
}
```

This example is also available at
GitHub (https://github.com/spark-
examples/spark-scala-
examples/blob/master/src/main/scala/co
m/sparkbyexamples/spark/dataframe/Gr
oupbyExample.scala) project for
reference.

## Conclusion

In this tutorial, you have learned how to
use `groupBy()` and aggregate
functions on Spark DataFrame and also
learned how to run these on multiple
columns and finally filtering data on the
aggregated column.

Thanks for reading. If you like it, please
do share the article by following the
below social links and any comments or
suggestions are welcome in the
comments sections!

Happy Learning !!

---

**Share this:**

**TAGS:** **AGG (HTTPS://SPARKBYEXAMPLES.COM/TAG/AGG/)**, **GROUPBY (HTTPS://SPARKBYEXAMPLES.COM/TAG/GROUPBY/)**

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

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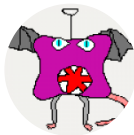
SparkByExamples.com is a Big Data and Spark examples community page, all examples are simple and easy to understand and well tested in our development environment Read more .. (https://sparkbyexamples.com/about-sparkbyexamples/)

### ❯ THIS POST HAS 3 COMMENTS

**divkr**    6 JAN 2021    REPLY

are you using python or scala for this tutorial ?

**NNK**    6 JAN 2021    REPLY

Examples on this page use Scala. If you are looking for GroupBy with Python (PySpark) see https://sparkbyexamples.com/pyspark/pyspark-groupby-explained-with-example/ (https://sparkbyexamples.com/pyspark/pyspark-groupby-explained-with-example/)

**sri bhargavi**

23 NOV 2020    REPLY

Hi, why we use agg without agg also can we perform agg functions rt??

## Leave a Reply

(https:/

/www.f    (https:/

(https:/    acebo    /www.li   (https:/

/twitter   ok.co    nkedin   /github

.com/s    m/spar   .com/i   .com/s

parkby   kbyexa   n/n-nk-   park-

examp   mples/   b860a   examp

les)        )       8193/)    les/)