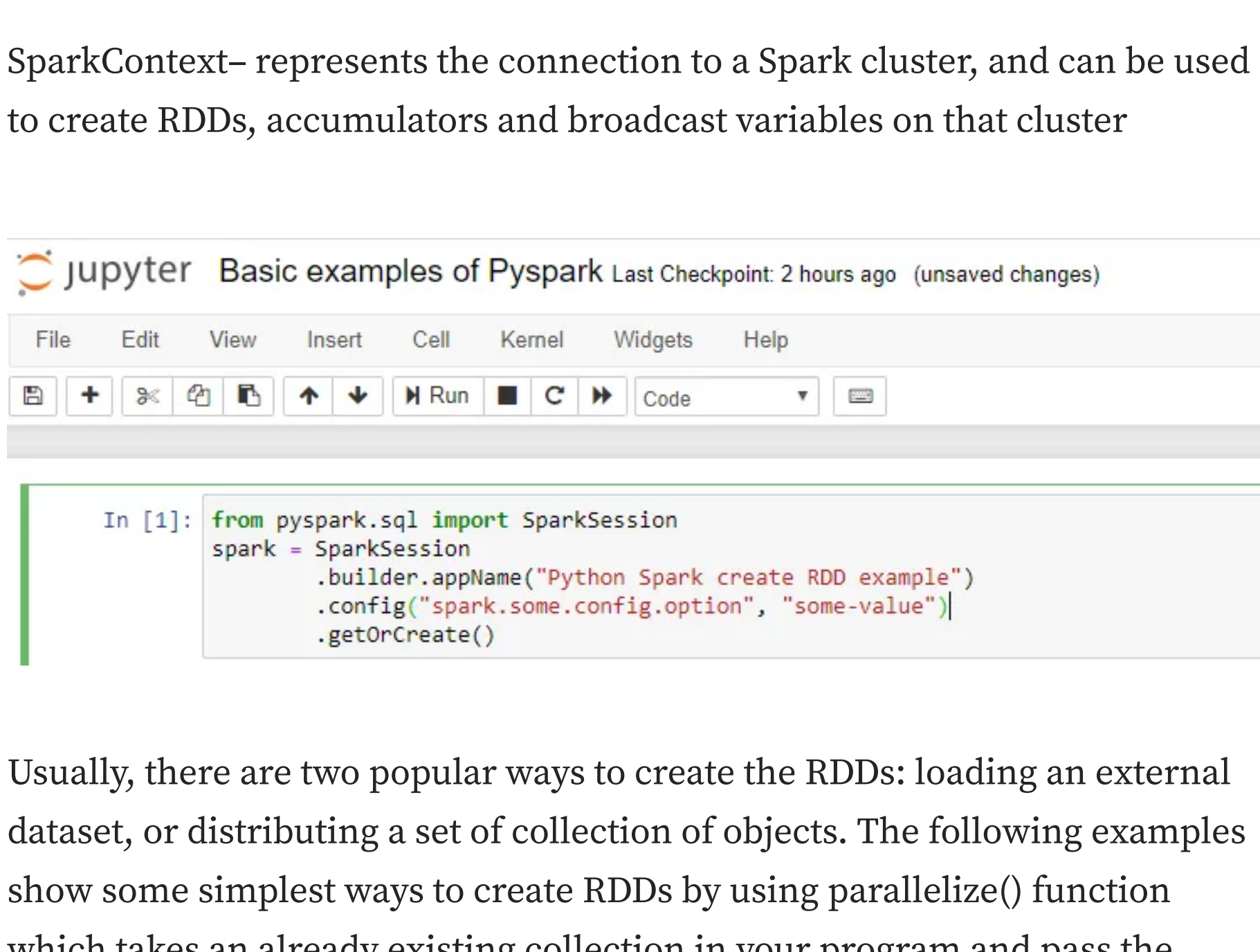


Basics of Pyspark Programming for RDD on Jupyter notebook

ISHMEET KAUR · Follow
3 min read · May 18, 2020



Usually, there are two popular ways to create the RDDs: loading an external dataset, or distributing a set of collection of objects. The following examples show some simplest ways to create RDDs by using parallelize() function which takes an already existing collection in your program and pass the same to the Spark Context. Then you will get RDD data.

```
myData = spark.sparkContext.parallelize([(1,2), (3,4), (5,6), (7,8), (9,10)])

myData.collect()

[1, 2], (3, 4), (5, 6), (7, 8), (9, 10)]
```

a. To Create Dataframe of RDD dataset:

1. With the help of toDF() function in parallelize function.

show() will display in the form of dataframe

collect() will display RDD in the list form for each row

```
In [15]: df = spark.sparkContext.parallelize([(1, 2, 3, 'a b c'), (4, 5, 6, 'd e f'), (7, 8, 9, 'g h i')])
         df.show()
         df.collect()

Out[15]: [Row(col1=1, col2=2, col3=3, col4=a b c),
          Row(col1=4, col2=5, col3=6, col4=d e f),
          Row(col1=7, col2=8, col3=9, col4=g h i)]
```

2. With createDataFrame() implicit call both arguments: RDD dataset can be represented in structured dataset with proper schema declared in the second argument of createDataFrame() of spark session.

```
Employee = spark.createDataFrame([
    ('1', 'Joe', '70000', '1'),
    ('2', 'Henry', '60000', '2'),
    ('3', 'Sam', '60000', '2'),
    ('4', 'Max', '90000', '1')],
    ['id', 'Name', 'Salary', 'DepartmentId'])

Employee.show()
```

Collect() will show RDD Row format.

```
Employee.collect()

[Row(id=1, Name=Joe, Salary=70000, DepartmentId=1),
 Row(id=2, Name=Henry, Salary=60000, DepartmentId=2),
 Row(id=3, Name=Sam, Salary=60000, DepartmentId=2),
 Row(id=4, Name=Max, Salary=90000, DepartmentId=1)]
```

3. With createDataFrame() explicitly call both the arguments: my_list having the list of values which will be set as rows and col_name is carrying the list of column names values.

```
: my_list = [['a', 1, 2], ['b', 2, 3], ['c', 3, 4]]
: col_name = ['A', 'B', 'C']

: spark.createDataFrame(my_list, col_name).show()

+----+----+
| A | B | C |
+----+----+
| a | 1 | 2 |
| b | 2 | 3 |
| c | 3 | 4 |
+----+----+
```

4. With createDataFrame() explicitly calls one of the argument: my_list dataset is called explicitly and column names are declared implicitly.

head()-displays first row of the RDD dataset.

show()-displays the dataframe

```
: my_list = [['male', 1, None], ['female', 2, 3], ['male', 3, 4]]

: ds = spark.createDataFrame(my_list, ['A', 'B', 'C'])
: ds.head()

: Row(A='male', B=1, C=None)

: ds.show()

+----+----+
| A | B | C |
+----+----+
| male | 1 | null |
| female | 2 | 3 |
| male | 3 | 4 |
+----+----+
```

5. createDataFrame() with numpy() array example-When dataset is in keys and values format.

```
d = {'A': [0, 1, 0],
     'B': [1, 0, 1],
     'C': [1, 0, 0]}

import numpy as np

spark.createDataFrame(np.array(list(d.values())).T.tolist(), list(d.keys())).show()

+----+----+
| A | B | C |
+----+----+
| 0 | 1 | 1 |
| 1 | 0 | 0 |
| 0 | 1 | 0 |
+----+----+
```

- b. fillna(value)-replace null value to the value provided in fillna() argument.

```
ds.fillna(-99).show()

+----+----+
| A | B | C |
+----+----+
| male | 1 | -99 |
| female | 2 | 3 |
| male | 3 | 4 |
+----+----+
```

- c. replace(listOfRowValues, listOfRowvaluesToReplace)-replacing row listed values with the list of values to replace with.

```
#caution: Mixed type replacements are not supported
ds.replace(['male', 'female'], ['1', '0']).show()

+----+----+
| A | B | C |
+----+----+
| 1 | 1 | null |
| 0 | 2 | 3 |
| 1 | 3 | 4 |
+----+----+
```

- d. toDF(listOfColvaluesToReplace)-To replace all the column names of dataframe with the new set of column names in the dataframes

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```

```
ds.toDF('id', 'name', 'salary', 'dept').show()

+----+----+
| id | name | salary | dept |
+----+----+
| 1 | Joe | 70000 | 1 |
| 2 | Henry | 60000 | 2 |
| 3 | Sam | 60000 | 2 |
| 4 | Max | 90000 | 1 |
+----+----+
```