



(<https://cognitiveclass.ai>)

Welcome to exercise one of “Apache Spark for Scalable Machine Learning on BigData”. In this exercise you’ll apply the basics of functional and parallel programming.

Let's start with a simple example. Let's consider you have a list of integers.

Let's find out what the size of this list is.

Note that we already provide an RDD object, so please have a look at the RDD API in order to find out what function to use: <https://spark.apache.org/docs/latest/api/python/pyspark.html#pyspark.RDD>

https://spark.apache.org/docs/latest/api/python/pyspark.html#pyspark.RDD?cm_mmc=Email_Newsletter_-_Developer_Ed%2BTech_-_WW_WW_-_SkillsNetwork-Courses-IBMDDeveloperSkillsNetwork-ML0201EN-SkillsNetwork-20647446&cm_mmca1=000026UJ&cm_mmca2=10006555&cm_mmca3=M12345678&cvosrc=email.Newsletter_-_Developer_Ed%2BTech_-_WW_WW_-_SkillsNetwork-Courses-IBMDDeveloperSkillsNetwork-ML0201EN-SkillsNetwork-20647446

The following link contains additional documentation: <https://spark.apache.org/docs/latest/rdd-programming-guide.html> (https://spark.apache.org/docs/latest/rdd-programming-guide.html?cm_mmc=Email_Newsletter-Developer_Ed%2BTech-WW-WW-SkillsNetwork-Courses-IBMDeveloperSkillsNetwork-ML0201EN-SkillsNetwork-20647446&cm_mmca1=000026UJ&cm_mmca2=10006555&cm_mmca3=M12345678&cvo_src=email.Newsletter)

This notebook is designed to run in a IBM Watson Studio default runtime (NOT the Watson Studio Apache Spark Runtime as the default runtime with 1 vCPU is free of charge). Therefore, we install Apache Spark in local mode for test purposes only. Please don't use it in production.

In case you are facing issues, please read the following two documents first:

<https://github.com/IBM/skillsnetwork/wiki/Environment-Setup>
<https://github.com/IBM/skillsnetwork/wiki/Environment-Setup>

<https://github.com/IBM/skillsnetwork/wiki/FAQ> (<https://github.com/IBM/skillsnetwork/wiki/FAQ>)

Then, please feel free to ask:

<https://coursera.org/learn/machine-learning-big-data-apache-spark/discussions/all>
https://coursera.org/learn/machine-learning-big-data-apache-spark/discussions/all?cm_mmc=Email_Newsletter-_Developer_Ed%2BTech-_WW_WW-_SkillsNetwork-Courses-IBMDeveloperSkillsNetwork-ML0201EN-SkillsNetwork-20647446&cm_mmca1=000026UJ&cm_mmca2=10006555&cm_mmca3=M12345678&cvosrc=email.Newsletter

Please make sure to follow the guidelines before asking a question:

<https://github.com/IBM/skillsnetwork/wiki/FAQ#im-feeling-lost-and-confused-please-help-me>
<https://github.com/IBM/skillsnetwork/wiki/FAQ#im-feeling-lost-and-confused-please-help-me>

If running outside Watson Studio, this should work as well. In case you are running in an Apache Spark context outside Watson Studio, please remove the Apache Spark setup in the first notebook cells.

In [1]:

```
from IPython.display import Markdown, display
def printmd(string):
    display(Markdown('# <span style="color:red">'+string+'</span>'))

if ('sc' in locals() or 'sc' in globals()):
    printmd('<<<<<!!!! It seems that you are running in a IBM Watson Studio Apache Spark Notebook. Please run it in an IBM Watson Studio Default Runtime (without Apache Spark) !!!!!>>>>>')
```

In [2]:

```
!pip install pyspark==2.4.5
```

Collecting pyspark==2.4.5

Downloading <https://files.pythonhosted.org/packages/9a/5a/271c416c1c2185b6cb0151b29a91fff6fcaed80173c8584ff6d20e46b465/pyspark-2.4.5.tar.gz> (217.8 MB)

|████████████████████████████████████████| 217.8MB 160kB/s eta 0:00:01

Collecting py4j==0.10.7 (from pyspark==2.4.5)

Downloading <https://files.pythonhosted.org/packages/e3/53/c737818eb9a7dc32a7cd4f1396e787bd94200c3997c72c1dbe028587bd76/py4j-0.10.7-py2.py3-none-any.whl> (197kB)

|████████████████████████████████████████| 204kB 48.0MB/s eta 0:00:01

Building wheels for collected packages: pyspark

Building wheel for pyspark (setup.py) ... done

Stored in directory: /home/dsxuser/.cache/pip/wheels/bf/db/04/61d66a5939364e756eb1c1be4ec5bdce6e04047fc7929a3c3c

Successfully built pyspark

Installing collected packages: py4j, pyspark

Successfully installed py4j-0.10.7 pyspark-2.4.5

In [3]:

```
try:
    from pyspark import SparkContext, SparkConf
    from pyspark.sql import SparkSession
except ImportError as e:
    printmd('<<<<<!!!! Please restart your kernel after installing Apache Spark !!!!!>>>>')
```

In [4]:

```
sc = SparkContext.getOrCreate(SparkConf().setMaster("local[*]"))

spark = SparkSession \
    .builder \
    .getOrCreate()
```

In [5]:

```
rdd = sc.parallelize(range(100))
```

In [6]:

```
# please replace $$ with the correct characters
rdd.count()
```

Out[6]:

100

You should see "100" as answer. Now we want to know the sum of all elements. Please again, have a look at the API documentation and complete the code below in order to get the sum.

In [7]:

```
rdd.sum()
```

Out[7]:

4950

You should get "4950" as answer.

Thank you for completing this lab!

This notebook was created by [Romeo Kienzler \(https://linkedin.com/in/romeo-kienzler-089b4557\)](https://linkedin.com/in/romeo-kienzler-089b4557) I hope you found this lab interesting and educational. Feel free to contact me if you have any questions!

Change Log

Date (YYYY-MM-DD)	Version	Changed By	Change Description
2020-09-29	2.0	Srishti	Migrated Lab to Markdown and added to course repo in GitLab

© IBM Corporation 2020. All rights reserved.

In []: