

MSAN 694 : Distributed Computing

Diane Woodbridge, Ph.D.
MSAN, University of San Francisco



Reviews

Class Overview

Motivation - Why Distributed Computing?

What is Distributed Computing?

Spark

RDD Creation

Spark Interview Questions

~~What is Apache Spark?~~

~~Explain the key features of Spark.~~

~~What is RDD?~~

~~How to create RDD.~~

~~What is "partitions"?~~

~~Types or RDD operations?~~

~~What is "transformation"?~~

~~What is "action"?~~

~~Functions of "spark core"?~~

~~What is "spark context"?~~

~~What is an "RDD lineage"?~~

~~Which file systems does Spark support?~~

~~List the various types of "Cluster Managers" in Spark.~~

~~What is "YARN"?~~

~~What is "Mesos"?~~

~~What is a "worker node"?~~

~~What is an "accumulator"?~~

~~What is "Spark SQL" (Shark)?~~

~~What is "SparkStreaming"?~~

~~What is "GraphX"?~~

~~What is "MLlib"?~~

<https://www.edureka.co/blog/interview-questions/top-apache-spark-interview-questions-2016/>

Spark Interview Questions

~~What are the advantages of using Apache Spark over Hadoop MapReduce for big data processing?~~

~~What are the languages supported by Apache Spark for developing big data applications?~~

~~Can you use Spark to access and analyze data stored in Cassandra databases?~~

Is it possible to run Apache Spark on Apache Mesos?

How can you minimize data transfers when working with Spark?

Why is there a need for broadcast variables?

~~Name a few companies that use Apache Spark in production.~~

What are the various data sources available in SparkSQL?

What is the advantage of a Parquet file?

What do you understand by Pair RDD?

Is Apache Spark a good fit for Reinforcement learning?

<https://www.dezyre.com/article/top-50-spark-interview-questions-and-answers-for-2016/208>

Contents

RDD Creation

RDD Operations - Transformations

RDD Operations - Actions

Run python script in Spark

Contents

RDD Creation

RDD Operations - Transformations

RDD Operations - Actions

Run python script in Spark

RDD Creation

Two ways of creating RDDs.

1. Loading an external data.

```
lines = sc.textFile("README.md")
```

2. Takes a collection such as `Seq` (`Array` or `List`) and creates RDD from its element and distribute to Spark executors in the process.

```
lines = sc.parallelize(["spark",  
"spark is fun!"])
```

- Check the number of partitions.

```
lines.getNumPartitions()
```

<https://spark.apache.org/docs/1.2.0/configuration.html>

Note :

Python Lambda Expression

Python Lambda Expression

- A shortened way to define functions inline.
- Create an anonymous function using the “lambda” keyword at runtime.

```
def f(x):  
    return x+2
```

f(2)

```
g = lambda x : x+2
```

g(2)

→ Can be used as a function parameter for RDD operations.

Contents

RDD Creation

RDD Operations - Transformations

RDD Operations - Actions

Run python script in Spark

RDD Operations

Two types

1. Transformation

- Produce a new RDD by performing data manipulation on another RDD.
- Ex. map, filter, flatmap, mapPartitions, sample, union, intersection, distinct, groupByKey, reduceByKey, aggregateByKey, sortByKey, join, cogroup, cartesian, pipe, coalesce, repartition, repartitionAndSortWithinPartitions.

2. Actions

- Trigger a computation to return the result to the calling program or to perform some actions on an RDD's elements.
- Ex. reduce, collect, count, first, take, takeSample, takeOrdered, saveAsTextFile, saveAsSequenceFile, saveAsObjectFile, countByValue, foreach.

<https://spark.apache.org/docs/1.2.0/programming-guide.html>

Contents

RDD Creation

RDD Operations - Transformations

RDD Operations - Actions

Run python script in Spark

RDD Operations - Transformation

Construct a new RDD from an existing RDD.

```
line_with_spark =  
lines.filter(lambda lines :  
"spark" in lines)
```

Lazy Evaluation. (cf. Hadoop MapReduce)

- Computation doesn't take place until an action is triggered.

Return RDDs.

RDD Operation - Transformation

Transformation Operation Types

<code>map(func)</code>	Return a new distributed dataset formed by passing each element of the source through a function <i>func</i> .
<code>filter(func)</code>	Return a new dataset formed by selecting those elements of the source on which <i>func</i> returns true.
<code>flatMap(func)</code>	Similar to map, but each input item can be mapped to 0 or more output items (so <i>func</i> should return a Seq rather than a single item).
<code>mapPartitions(func)</code>	Similar to map, but runs separately on each partition (block) of the RDD, so <i>func</i> must be of type <code>Iterator<T> => Iterator<U></code> when running on an RDD of type T.
<code>mapPartitionsWithIndex(func)</code>	Similar to mapPartitions, but also provides <i>func</i> with an integer value representing the index of the partition, so <i>func</i> must be of type <code>(Int, Iterator<T>) => Iterator<U></code> when running on an RDD of type T.
<code>sample(withReplacement, fraction, seed)</code>	Sample a fraction <i>fraction</i> of the data, with or without replacement, using a given random number generator seed.
<code>distinct([numTasks])</code>	Return a new dataset that contains the distinct elements of the source dataset.
<code>union(otherDataset)</code>	Return a new dataset that contains the union of the elements in the source dataset and the argument.
<code>intersection(otherDataset)</code>	Return a new RDD that contains the intersection of elements in the source dataset and the argument.

RDD Operation - Transformation

Element-wise Transformation

- `map(func)`
 - Apply a function to each element in the RDD.
- `flatMap(func)`
 - Call each element in RDD individually.
 - Concatenates multiple arrays into a collections that has one level structure.
- `filter(func)`
 - Return an RDD that passes the filtering requirement.

Example 1-1

Load a text file(“ignatian_pedagogy”) and split each line by space.

Example 1-1

Load a text file("ignatian_pedagogy") and split each line by space.

```
lines = sc.textFile("../Data/ignatian_pedagogy")  
  
lines.collect()  
  
words = lines.map(lambda line : line.split())  
  
words.collect()
```

Create an RDD representing the lines of text in a file.

Example 1-2

Generate a list of words within one level structure.

Example 1-3

Find words including “USF”.

RDD Operation - Transformation

Partition-wise Transformation

- `mapPartitions(func)`
 - Return a new RDD by applying a function to each partition of the RDD.
- `mapPartitionsWithIndex(func)`
 - Return a new RDD by applying a function to each partition of the RDD, while tracking the index of the original partition.

Example 2

Parallelize numbers between 1 and 16.

Calculate the count and sum in each partition.

RDD Operation - Transformation

Set Operation

- Format : `rdd1.operator(rdd2)`
- `distinct()`
 - Return only one of each element.
- `union()`
 - If there are duplicated elements, it returns all duplicates.
- `intersection()`
 - Return common elements.
- `subtract()`
 - Return elements that are in `rdd1` only.
- `cartesian()`
 - Return cartesian product (all pairs between `rdd1` and `rdd2`)

Example 3-1

Find distinct words in “ignatian_pedagogy”.

Example 3-2

Create a flatmap of distinct words from
“README.md”

Example 3-3

What is union, intersection, subtract and cartesian product of the sets from Example 3-1 and Example 3-2?

Contents

RDD Creation

RDD Operations - Transformations

RDD Operations - Actions

Run python script in Spark

RDD Operation - Action

Compute a result based on an RDD.

Return the result to the driver program or save it to external storage system.

```
line_with_spark.count()
```

Return non-RDDs.

RDD Operation - Action

Action Operation Types

<code>reduce(func)</code>	Combine the elements of the RDD together in parallel. (eg. Sum)
<code>fold(zero)(func)</code>	Same as <code>reduce()</code> , but with the provided zero value.
<code>aggregate(zero)(SeqOp, combOp)</code>	Similar to <code>reduce()</code> but used to return a different type.
<code>collect()</code>	Return all the elements of the dataset as an array at the driver program. This is usually useful after a filter or other operation that returns a sufficiently small subset of the data.
<code>count()</code>	Return the number of elements in the dataset.
<code>countByValue()</code>	Return the number of times each element occurs in the RDD.
<code>take(n)</code>	Return an array with the first <code>n</code> elements of the dataset.
<code>top(n)</code>	Return the top <code>n</code> elements of the RDD.
<code>first()</code>	Return the first element of the dataset (similar to <code>take(1)</code>).
<code>takeSample(withReplacement, num, [seed])</code>	Return an array with a random sample of <code>num</code> elements of the dataset, with or without replacement, optionally pre-specifying a random number generator seed.
<code>foreach(func)</code>	Run a function <code>func</code> on each element of the dataset. This is usually done for side effects such as updating an Accumulator or interacting with external storage systems. Note: modifying variables other than Accumulators outside of the <code>foreach()</code> may result in undefined behavior. See Understanding closures for more details.

RDD Operation - Action

- `reduce(func)`
 - Take a function that operates on two elements of the type in your RDD.
 - Returns a new element of the same type.
Ex. `sum = rdd.reduce(lambda x, y = x+y)`
- `fold(zeroValue)(func)`
 - Take a function with the same signature as needed for `reduce()`.
 - Take a “zero value” to be used for the initial call on each partition.
 - Returns a new element of the same type.
- `aggregate(zeroValue)(seqOp, combOp)`
 - Supply an initial zero value of the type we want to return.
 - `seqOp` : Function to combine the elements from the RDD with the accumulator. Runs once in a partition.
 - `combOp` : Function to merge two accumulators, given that each nodes accumulates its own results locally.
 - Can return an element of a different type.

Example 4-1

For the numbers between 1 and 9, calculate sum of the odd numbers.

Example 4-2

For the numbers between 1 and 9, calculate sum of the odd numbers using fold().

Example 4-3

Using `aggregate()`, return (sum, # of elements) of odd numbers.

RDD Operation - Action

- `collect()`
 - Return the entire RDD's contents.
- `count()`
 - Return the count of elements.
- `countByValue()`
 - Return the number of times each element occurs in the RDD.
- `top(n)`
 - Return top elements from an RDD, using the default ordering on the data.
- `take(n)`
 - Return n elements.
- `first()`
 - Return the first element of the data.
- `takeSample(withReplacement, num, seed)`
 - Return a fixed-size sample subset of an RDD.
- `foreach()`
 - Used for performing computations on each element in the RDD.

RDD Operation - Action

sample() vs. takeSample()

- sample(withReplacement, fraction, seed) : **Transformation**
 - Creates a new RDD with random elements from the calling RDD.
 - withReplacement : Allow sample multiple times.
 - fraction :
 - Expected number of times each element is going to be sampled (positive double), when replacement is used.
 - Expected probability that each element is going to be sampled (between 0 and 1), when replacement is not used.
 - seed : Random number generation. (Same seeds generates the same numbers.)

RDD Operation - Action

sample() vs. takeSample()

- takeSample(withReplacement, num, seed) : **Action**
 - Return a fixed-size sample subset of an RDD as an array (not RDD).
 - withReplacement : Allow sample multiple times.
 - num : The exact number of sampled element. (Integer)
 - seed : Random number generation. (Same seeds generates the same numbers.)

Example 5-1

```
x = sc.parallelize([3,4,1,2])  
y = sc.parallelize(range(2,6))  
z = x.union(y)
```

Try collect(), count(), countByValue(), top(n), take(n), first(), takeSample() operations on z.

RDD Operation - Action

Numeric RDD Action Types

<code>count()</code>	Return the number of elements in the RDD.
<code>mean()</code>	Return the mean of the RDD's elements.
<code>sum()</code>	Add up the elements in the RDD.
<code>max()</code>	Return the maximum item in the RDD.
<code>min()</code>	Return the minimum item in the RDD.
<code>variance()</code>	Return the variance of the RDD's elements.
<code>stdev()</code>	Return the standard deviation of the RDD's

<https://spark.apache.org/docs/1.1.1/api/python/pyspark.rdd.RDD-class.html>

Contents

RDD Creation

RDD Operations - Transformations

RDD Operations - Actions

Run python script in Spark

Run python script in Spark

```
from pyspark import SparkConf, SparkContext
```

Create SparkConf object to configure the application.

```
conf = SparkConf().setMaster("local[*]").setAppName("AppName")
```

Initializing a SparkContext (SC).

```
sc = SparkContext(conf = conf)
```

For closing Spark, call `sc.stop()`

Unset your environment variables.

Run your standalone program.

```
spark-submit yoursript
```

Run python script in Spark

```
from pyspark import SparkConf, SparkContext
```

Create SparkConf object to configure the application.

```
conf = SparkConf().setMaster("local[*]").setAppName("AppName")
```

Initializing a SparkContext (SC).

```
sc = SparkContext(conf = conf)
```

Cluster URL

"local[*]" : indicate all thread on the local machine.

For closing Spark, call `sc.stop()`

Unset your environment variables.

Run your standalone program.

```
spark-submit yoursript
```

Run python script in Spark

```
from pyspark import SparkConf, SparkContext
```

Create SparkConf object to configure the application.

```
conf = SparkConf().setMaster("local[*]").setAppName("AppName")
```

Initializing a SparkContext (SC).

```
sc = SparkContext(conf = conf)
```

For closing Spark, call `sc.stop()`

Unset your environment variables.

Run your standalone program.

```
spark-submit yoursript
```

Identify the application if you're connecting to a cluster.

Example 6

Write a python script (.py) for printing the number of lines in "README.md" and the first line and run on spark.

Example 6

```
from pyspark import SparkConf, SparkContext

#Create SparkContext
conf =
SparkConf().setMaster("local[*]").setAppName("Ap
pName")

sc = SparkContext(conf = conf)

#Load Data.
lines=sc.textFile("../Data/README.md")
print(lines.count())
print(lines.first())

sc.stop()
```

```
unset PYSPARK_DRIVER_PYTHON
unset PYSPARK_DRIVER_PYTHON_OPTS
spark-submit ex6.py > output_file.txt
```

References

Distributed Computing with Spark, Reza Zadeh,
http://stanford.edu/~rezab/slides/bayacm_spark.pdf

Spark Online Documentation : [http://
spark.apache.org/docs/latest/](http://spark.apache.org/docs/latest/)

Spark RDD Class: [https://spark.apache.org/docs/
1.1.1/api/python/pyspark.rdd.RDD-class.html](https://spark.apache.org/docs/1.1.1/api/python/pyspark.rdd.RDD-class.html)

Karau, Holden, et al. Learning spark: lightning-fast big data analysis. O'Reilly Media, Inc., 2015.

Zecevic, Petar, et al. Spark in Action, Manning, 2016.