Learning the Core of Spark: RDDs



Xavier Morera
HELPING DEVELOPERS UNDERSTAND SEARCH & BIG DATA
@xmorera www.xaviermorera.com



Learning the Core of Spark: RDDs



Resilient Distributed Dataset

How can I use DataFrames in 2.0

What is an RDD and Schema RDD

How do I group by a field

Can I use Hive from HUE





Why Start with RDDs?

Knowing RDDs gives you a better foundation for writing Spark Code

Sometimes RDDs are the right tool for the job

Existing Spark 1.x code



SparkContext Driver program Context Job Cluster Manager Worker Node Worker Node Worker Node Executor Executor Executor Task Task Task Task Task

You will need it to work with RDDs!



Spark Application

- Created for you in the shell
- You need to create in an application

Only one SparkContext





Spark Application



Configuration & Environment



Services



Create RDDs



```
sc
dir(sc)
sc.serializer
sc.sparkUser
help(sc) sc.sparkUser
sc.sparkUser()
sc.appName
sc.applicationId
```

SparkContext available as sc

Many methods and properties available

Get configuration and environment information



```
sc.stop()
test_rdd = sc.emptyRDD()
sc =
SparkContext.getOrCreate()
```

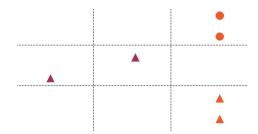
Stop and create a new SparkContext

- getOrCreate()

Used to create RDDs



RDD



How can I use <u>DataFrames</u> in 2.0

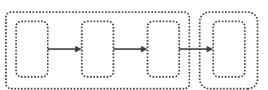
What is an RDD and Schema RDD

How do I group by a field

Can I use Hive from HUE

How can I use DataFrames in 2.0
What is an RDD and Schema RDD
How do I group by a field
Can I use Hive from HUE

Resilient

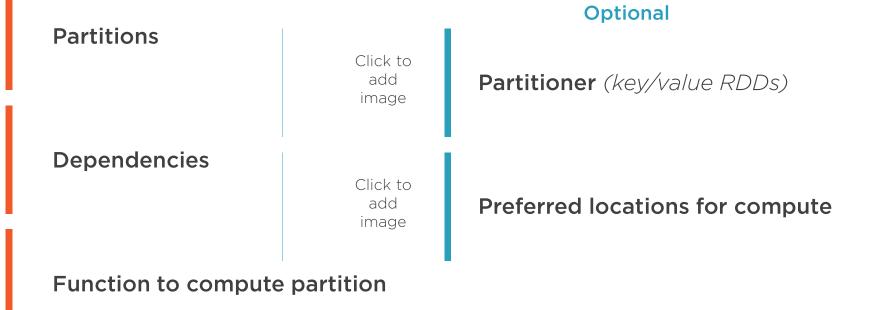


Distributed

Dataset



Five Properties of an RDD



```
* Internally, each RDD is characterized by five main properties:
62
63
64
     * - A list of partitions
     * - A function for computing each split
65
     * - A list of dependencies on other RDDs
66
     * - Optionally, a Partitioner for key-value RDDs (e.g. to say that the RDD is hash-partitioned)
67
     * - Optionally, a list of preferred locations to compute each split on (e.g. block locations for
68
69
          an HDFS file)
70
     * All of the scheduling and execution in Spark is done based on these methods, allowing each RDD
71
     * to implement its own way of computing itself. Indeed, users can implement custom RDDs (e.g. for
72
73
     * reading data from a new storage system) by overriding these functions. Please refer to the
74
     * <a href="http://people.csail.mit.edu/matei/papers/2012/nsdi_spark.pdf">Spark paper</a>
     * for more details on RDD internals.
75
76
     */
77
    abstract class RDD[T: ClassTag](
        @transient private var _sc: SparkContext,
78
79
        @transient private var deps: Seg[Dependency[_]]
      ) extends Serializable with Logging {
80
81
      if (classOf[RDD[_]].isAssignableFrom(elementClassTag.runtimeClass)) {
82
83
        // This is a warning instead of an exception in order to avoid breaking user programs that
        // might have defined nested RDDs without running jobs with them.
84
        logWarning("Spark does not support nested RDDs (see SPARK-5063)")
      }
```

PairRDD

RDDs of key/value pairs

Contain tuples

Useful for grouping or aggregating

Can use RDD transformations

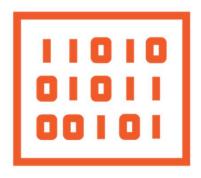
PairRDD specific transformations available



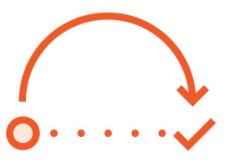
Creating RDDS



Parallelize



External Data



* or DataFrame (next module)



```
[1,2,3,4,5]
type[1,2,3,4,5]
list_one_to_five =sc.parallelize([1,2,3,4,5])
list_one_to_five
list_one_to_five.collect()
```

Creating RDDs with Parallelize

Create RDDs with data in memory: collection or iterable

Parallelize a collection

Bring data back to driver with collect()



```
list_one_to_five.getNumPartitions()
list_one_to_five.glom().collect()
list_one_to_five = sc.parallelize([1,2,3,4,5], 1)
list_one_to_five.getNumPartitions()
list_one_to_five.glom().collect()
```

RDD Partitions with Parallelize

Call getNumPartitions() to check number of partitions

Use glom() to see our data within partitions

And you can control how many partitions you want



dir(list_one_to_five)
list(list_one_to_five)

RDD Functionality

Use dir() to inspect attributes

Even better, use help()



```
list_one_to_five.sum()
list_one_to_five.min()
list_one_to_five.max()
list_one_to_five.mean()
```

An Operation on an RDD

What can I do on an RDD?

Call sum() to add up values

And first() to get first element



```
list_dif_types = sc.parallelize([False, 1, "two", {'three': 3},
    ('xavier', 4)])
list_dif_types.collect()
tuple_rdd = sc.parallelize([('xavier', 1),('irene', 2)])
tuple_rdd.setName('tuple_rdd')
tuple_rdd
```

Different Types of Objects

RDDs can hold different types of objects

Strings, ints, dictionaries, and key/value pairs (tuple) are common

RDD is given an internal name, but you can set it too



```
empty_rdd=sc.parallelize([])
empty_rdd.isEmpty()
empty_rdd=sc.emptyRDD()
empty_rdd.isEmpty()
```

EmptyRDD

You can create empty RDDs

With parallelize()

Or emptyRDD()



```
bigger_rdd = sc.parallelize(range(1,1000))
bigger_rdd.collect()
bigger_rdd.sum()
bigger_rdd.count()
bigger_rdd.getNumPartitions()
```

Ranges

Create RDDs with range()

And use methods to get information from our RDD data

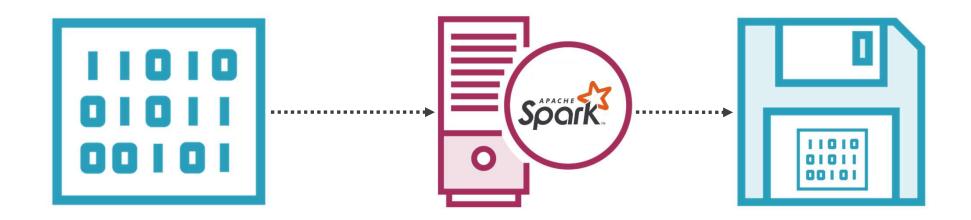
Count() all elements in my rdd

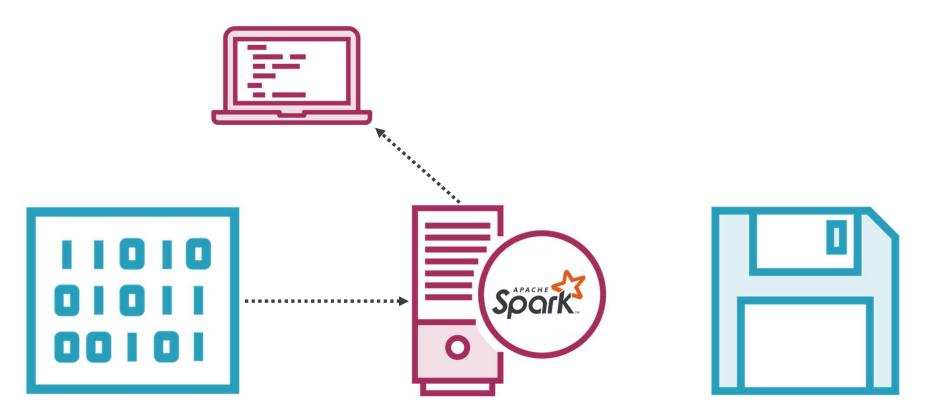


Parallelize

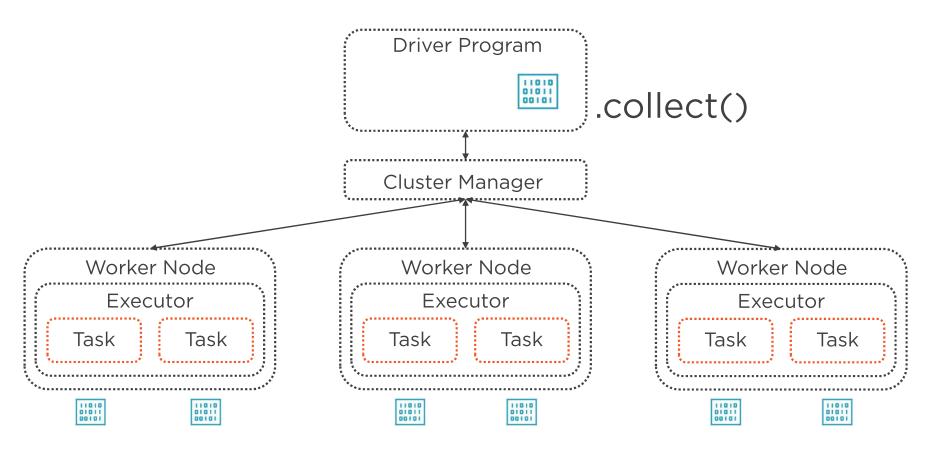












```
list_bigger.collect()
list_bigger.take(10)
list_bigger.first()
help(list_bigger.takeOrdered)
list_bigger.takeOrdered(10, key=lambda x: -x)
tuple_map = tuple_rdd.collectAsMap()
type(tuple_map)
tuple_map['xavier']
```

Bring data back RDD to Driver

Most common action for testing is collect()

Also first(), take(), takeOrdered(), collectAsMap() among others



ABeware!

Potential out of memory exception on large datasets



```
for elem in list_bigger.take(10):
    print elem

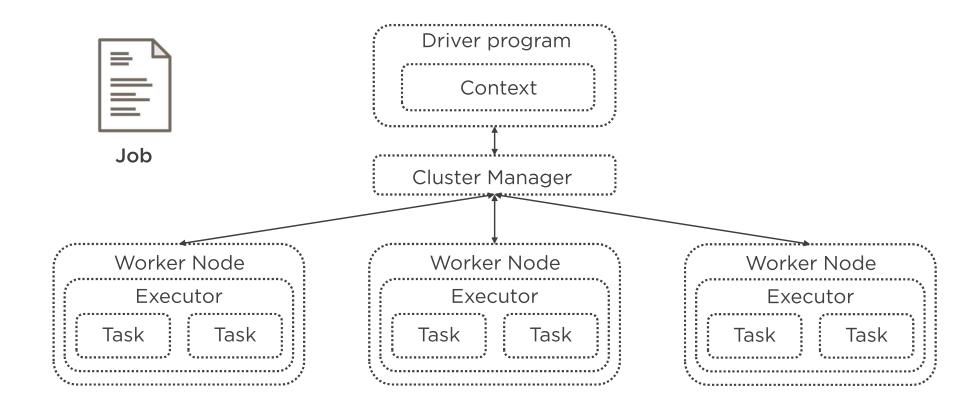
for key, value in tuple_map.iteritems():
    print key + " / " + str(value)
```

Iterating & Printing Data

Iterate over list of data returned using take() in for loop

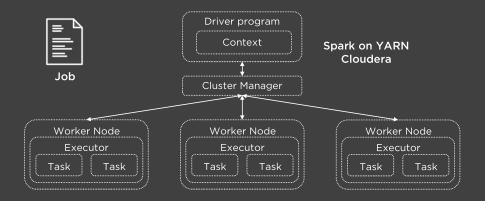
Or other types, i.e. dict





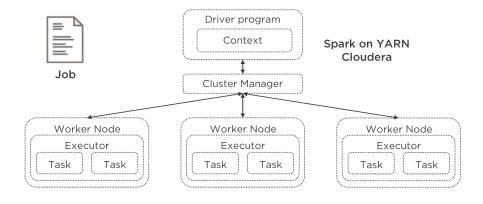
Driver

print test_write.collect()



In executor

```
def in_executor(entry):
    print "*** inExecutorEntry ***"
    print entry[0]
    print type(entry)
    return 'printed:'+str(entry[0])
```



spark2-submit inExecutor EntryLogLevelAll.py



```
import urllib2
def call_search(q):
    response = urllib2.urlopen('http://tiny.bigdatainc.org/' + q)
html = response.read
topics = sc.parallelize(['ts1', 'ts2'])
topics.foreach(call_search)
```

Foreach

Executes the function on each element in the RDD

Runs in the worker nodes

Useful for calling external resources like a database or an API



```
bigger_rdd.getNumPartitions()
play_part=bigger_rdd.repartition(10)
play_part.getNumPartitions()
```

(Re)Partitions, Coalesce, Saving Text & HUE bigger_rdd was partitioned on parallelize()

Change number of partitions with repartition() and coalesce()



```
play_part.repartition(14).getNumPartitions()
play_part.coalesce(15).getNumPartitions()
```

Repartition & Coalesce

repartition() specifies new number of partitions, up or down coalesce() uses existing partitions, only decreases, no shuffling Data in RDD remains the same. How can we tell?



```
play_part.saveAsTextFile("/user/cloudera/tests/
play_part/ten")

play_part.repartition(4).saveAsTextFile("/user/cloudera/
tests/play_part/four")

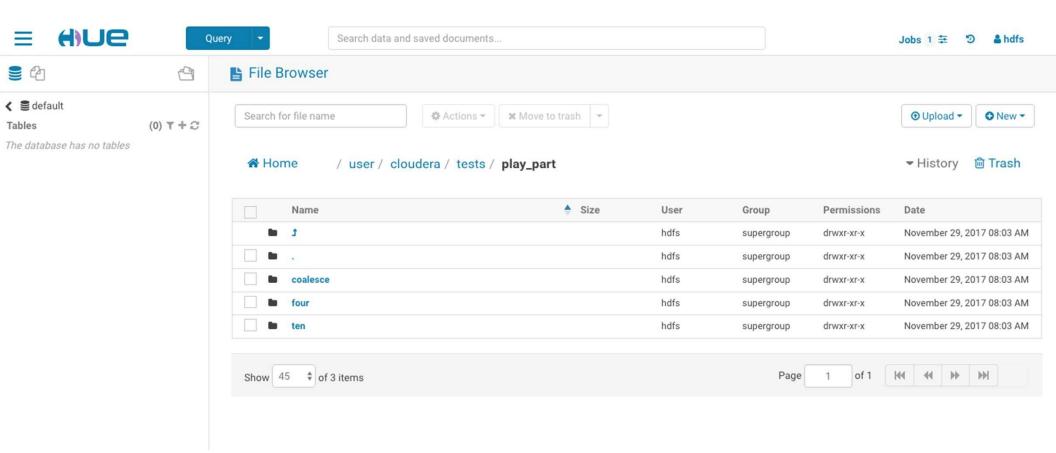
play_part.coalesce(1).saveAsTextFile("/user/cloudera/
tests/play_part/coalesce")
```

Confirm Data in RDD

Save using saveAsTextFile() and visualize with HUE

Just like Hadoop, creates file per partition







1. Warning!

Exception if folder already exists



play_part.coalesce(1).saveAsTextFile("/user/cloudera/tests/
play_part/coalesce")

Folder Exists

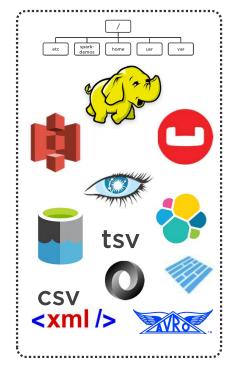
If folder exists, you get an exception

This applies for lower level API

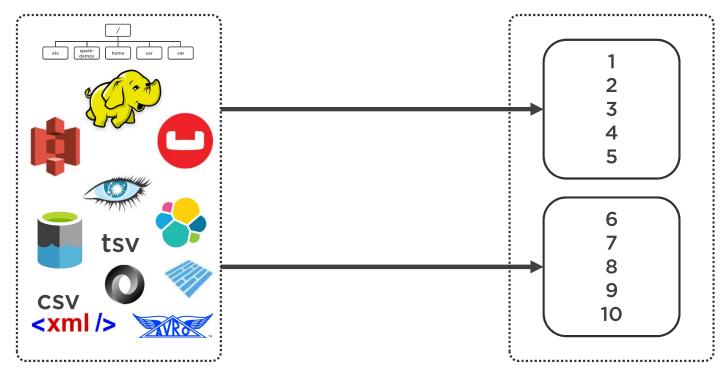
Works different for higher level API











```
sc.textFile("/user/cloudera/tests/play_part/coalesce")
.count()
sc.textFile("/user/cloudera/tests/play_part/ten").count()
sc.textFile("/user/cloudera/tests/play_part/four").count()
sc.textFile("/user/cloudera/tests/play_part/four/
part-00000").count()
```

Loading Text from HDFS

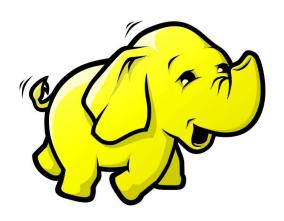
Use textFile(), using the SparkContext

Use hdfs://with.adifferent.cluster.manager

All files in the folder, or specific file



HDFS



Hadoop Distributed File System

Primary storage system for Hadoop

Data replication

Blocks

Optimized for parallel data processing



```
local_play =
sc.textFile("file:///stackexchange/play_part/")
local_play.count()
local_play.take(10)
```

Loading Local Data

Access local data using file:///

File needs to be in all nodes

Not required in standalone Spark



Convert Posts.xml to CSV



Data preparation step



```
posts_all =
sc.textFile("/user/cloudera/stackexchange/posts_all_csv")
posts_all.count()
```

Loading StackExchange / StackOverflow Posts
Load from HDFS

How many do we have?



CSV



Comma Separated Values

Plain text file

Tabular data

Supports different separators

- Tab



Convert Badges.xml to CSV



Data preparation step



Loading CSV

Loaded like any other file, with textFile()

You get lines, and tell Spark what to do with the data

Advantage and drawback at the same time



```
tags_partitions=sc.wholeTextFiles("/user/cloudera/stackexchange/badges_csv")
tags_partitions.take(1)
numbers_partitions = sc.wholeTextFiles("/user/cloudera/tests/play_part/four")
numbers_partitions.take(1)
numbers_partitions.take(1)[0][0]
numbers_partitions.take(1)[0][1]
```

Loading Whole Text Files

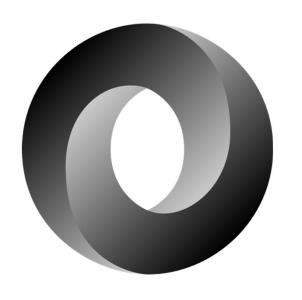
Loads each file into a row, uses wholeTextFiles()

In a PairRDD

Key is file name, value is the file



JSON



JavaScript Object Notation

Light weight data format

Easy to use

Collection of name value pairs and arrays

Fat-free XML



Convert Tags.xml to JSON



Data preparation step



```
tags_json=sc.textFile("/user/cloudera/stackexchange/
tags_json")
tags_json.first()
```

Load JSON

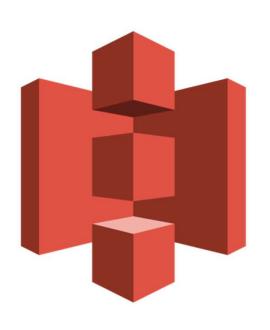
Just like CSV, it loads as text

Parse each record

Is there an easy way to parse JSON with RDDs?



S3



Simple Storage Service

Amazon's cloud storage

Buckets

Rock solid SLA

Unlimited scalability

At extremely reasonable cost

Works well with many technologies



Manually upload a JSON file into an S3 account
Get the full path in S3
Create an IAM user, enable programmatic access
Obtain access key id and secret access key
Set keys in environment
Load necessary packages

Prepare Data for Processing in S3



Data preparation step



```
# export AWS_ACCESS_KEY_ID="access-key"
# export AWS_SECRET_ACCESS_KEY="secret-key"
tags_json_s3 = sc.textFile("s3a://pscs/tags_json/part-00000.json")
tags_json_s3tags_json_s3.take(3)
```

Accessing Amazon S3 with Spark

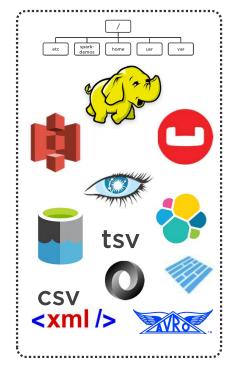
Transparent, still textFile()

Use s3a://

Be very careful with security







badges_columns_rdd.take(1)

badges_columns_rdd.saveAsTextFile("/user/cloudera/stackexchan
ge/badgestxt")

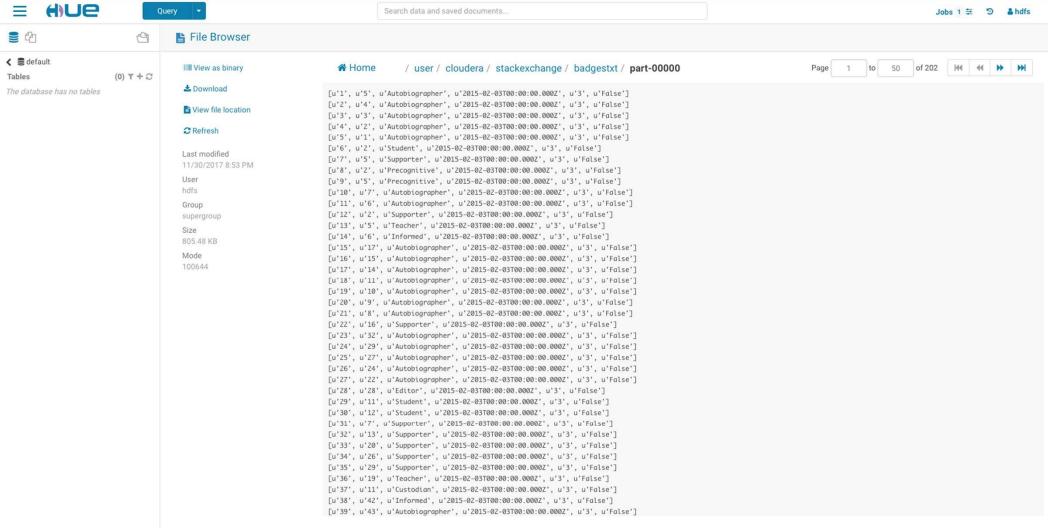
Saving an RDD

We previously saved an RDD that had no transformations applied

- Partition files

Save an RDD with columns as text







```
badges_cr_rdd =
sc.textFile("/user/cloudera/stackexchange/badgestxt")
badges_columns_rdd
badges_cr_rdd
```

Saving & Reloading an RDD

Reload our data

What happened when I save?



badges_columns_rdd.take(1)[0][2]
badges_cr_rdd.take(1)[0][2]

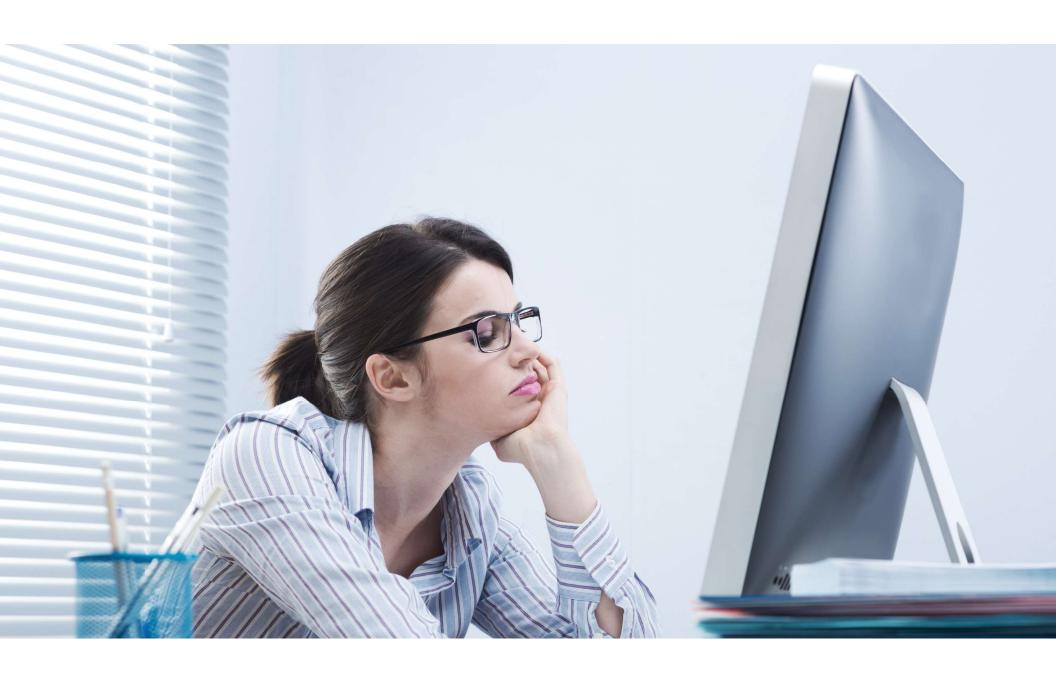
Text Is Good but Not Great

Compare original RDD vs. the one we just reloaded

We lost the format we applied

Text has no schema





Pickle in Python



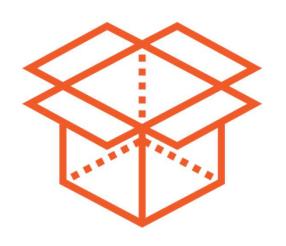
Serializing and deserializing objects

- Serializing → pickling
- Deserializing → unpickling

Preserve the object

Used regularly in Python

Pickle Files



Available in Spark too

Save as Pickle files

Do not lose applied transformations

```
badges_columns_rdd.saveAsPickleFile("/user/cloudera/stackex
change/badgespickle")

badges_crp_rdd =
sc.pickleFile("/user/cloudera/stackexchange/badgespickle")

badges_crp_rdd.take(1)

badges_crp_rdd.take(1)[0][2]
```

PickleFile

Save with saveAsPickleFile()

And load with pickleFile()



Saving a PickleFile in Spark actually saves a SequenceFile of Pickle objects

Watch out for this...



Sequence Files



Flat file

Binary format

Key/value pairs

Used extensively in Hadoop



```
badges_sorted.take(1)
badges_sorted.saveAsSequenceFile("/user/cloudera/
stackexchange/badgessequence")
badges_ss = sc.sequenceFile("/user/cloudera/stackexchange/
badgessequence")
badges_ss.take(1)
```

SequenceFile

Save using saveAsSequenceFile()

Read with sequenceFile()



Hadoop Formats



InputFormat

- How to read input files
- How they are split

OutputFormat

- How to write output files

Very flexible

```
badges_count_badge
   .saveAsNewAPIHadoopFile('/user/cloudera/badgesseqhadoop',
   "org.apache.hadoop.mapreduce.lib.output.SequenceFileOutputFormat",
   "org.apache.hadoop.io.IntWritable",
   "org.apache.hadoop.io.Text")
```

NewAPIHadopFile

More parameters required than other methods

i.e. file output format, key class, value class

And configuration



```
xml_posts =
   sc.newAPIHadoopFile('/user/cloudera/stackexchange/Posts.xml',
   'com.databricks.spark.xml.XmlInputFormat',
   'org.apache.hadoop.io.Text',
   'org.apache.hadoop.io.Text',
   conf={'xmlinput.start':'<row','xmlinput.end':'/>'})
```

Does newAPIHadoopFile Look Familiar?

We used to load the StackExchange dumps

Using spark-xml package





textFile

pickleFile

sequenceFile

hadoopFile

hadoopDataset

newAPIHadoopFile



saveAsTextFile
saveAsPickleFile
saveAsSequenceFile
saveAsHadoopFile
saveAsHadoopDataset
saveAsNewAPIHadoopFile





Full stop... nothing else to say...



```
rdd_reuse = sc.parallelize([1,2])
rdd_reuse.collect()
rdd_reuse = sc.parallelize([3,4])
rdd_reuse.collect()
```

But... (On Immutable)

Immutable

- Can't change the data

Reuse the variable



[u'1,5,Autobiographer,2015-02-03T00:00:00.000Z,3,False']



[u'1', u'5', u'Autobiographer', u'2015-02-03T00:00:00.000Z', u'3', u'False']



```
badges_rdd_csv.take(1)
split_the_line
badges_columns_rdd.take(1)
```

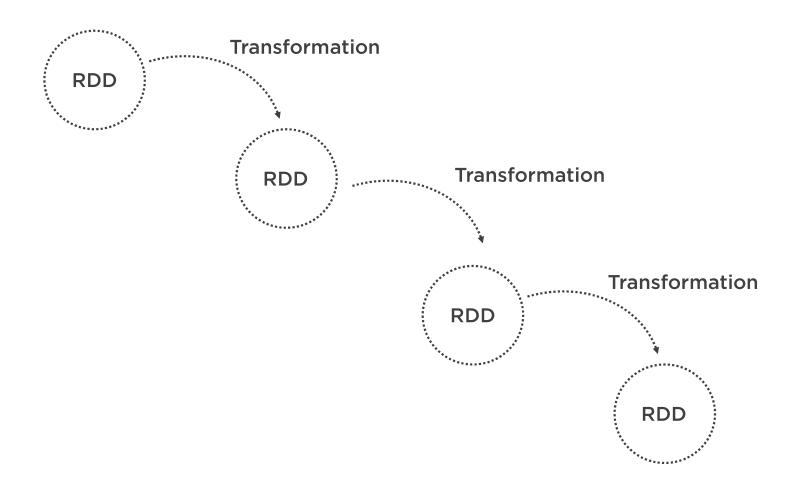
Creating RDDs with Transformations

Loaded Badges CSV into an RDD

- Single element per row

Applied map(), with split_the_line() function





```
badges_entry = badges_columns_rdd.map(lambda x: x[2])
badges_name = badges_entry.map(lambda x: (x,1))
badges_reduced = badges_name.reduceByKey(lambda x, y: x + y)
badges_count_badge = badges_reduced.map(lambda (x,y): (y,x))
badges_sorted = badges_count_badge.sortByKey(False).map(lambda (x, y): (y, x))
badges_sorted.take(10)
```

Creating Multiple RDDs with Transformations One transformation

Iterate

- Applying multiple transformations
- One of the main objectives of Spark



```
from pyspark import RDD
locals()
globals()
vars()['badges_entry']
vars()['badges_entry'].count()
for(k,v)inlocals().items():
   if isinstance(v,RDD) and 'badge' in k:
        print k
```

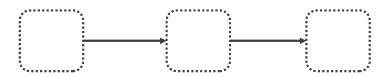
Tip: Listing RDDs

Check for created RDDs with globals(), locals() or vars()

Look for those that are an instance of RDD

- Use if

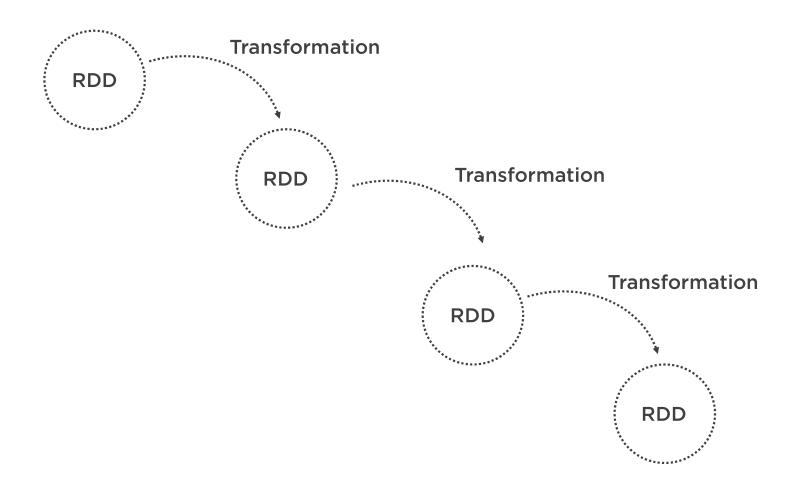




RDD Lineage

Graph of transformation operations required to execute when an action is called

RDD operator graph or RDD dependency graph

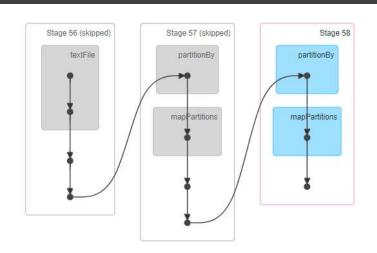


```
badges_sorted.toDebugString()
print badges_sorted.toDebugString()
badges_sorted
badges_sorted.prev
badges_sorted.take(1)
badges_sorted.prev.take(1)
```

A Little Bit More on Lineage Use toDebugString()

And print to get nice formatting

Just for fun, use prev





PySparkShell application UI

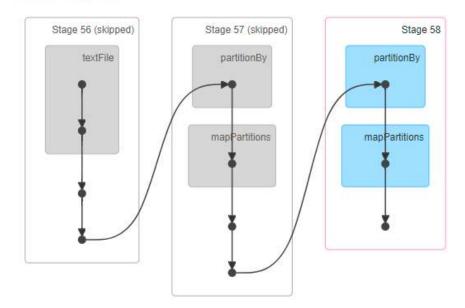


Details for Job 44

Status: SUCCEEDED Completed Stages: 1 Skipped Stages: 2

▶ Event Timeline

▼ DAG Visualization



Completed Stages (1)

| Stage Id * | Description | | Submitted | Duration | Tasks: Succeeded/Total | Input | Output | Shuffle Read | Shuffle Write |
|------------|-------------------------------|----------|---------------------|----------|------------------------|-------|--------|--------------|---------------|
| 58 | runJob at PythonRDD.scala:446 | +details | 2018/01/12 06:38:38 | 0.9 s | 1/1 | | | 803.0 B | |



Why learn RDDs?

- Better foundation
- Right tool for the job
- Existing code base



SparkContext

One context per application

RDD main abstraction

Multiple types of RDDs





Create RDDs

- Data in memory
- Files in storage
- Other RDDs

Different data sources and file formats

- Text, CSV, Parquet, Picklefile, JSON, Sequence Files among others
- Transparent





Return data back to the Driver

Persist to disk

- Some formats store schema
- Others don't

RDD Lineage

