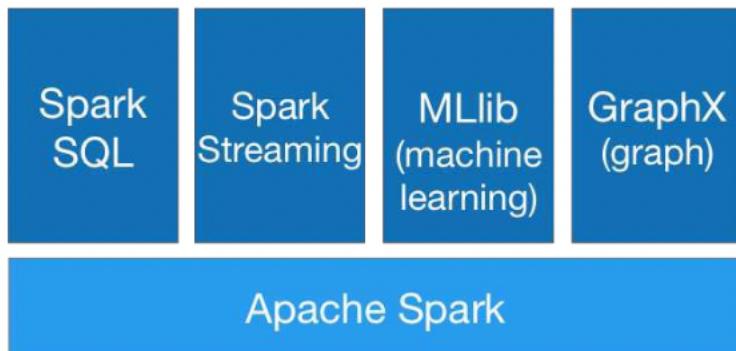


Spark Pair RDD Data Frame

Spark Architecture



Easy Ways to Run Spark

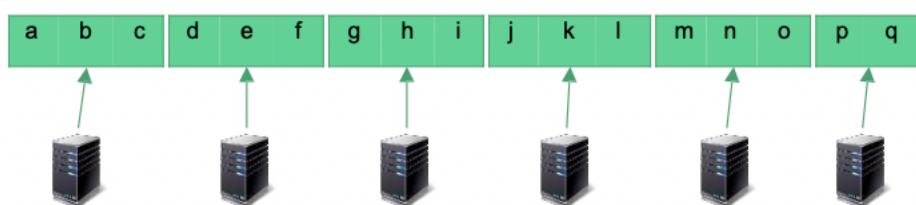
- IDE.
- Standalone Deploy Mode.
- Docker & Zeppelin.
- EMR.
- Hadoop vendors ([Cloudera](#), [Hortonworks](#)).
- Digital Ocean (Kubernetes cluster).

Supported Languages



RDD

- An RDD is simply an **immutable** distributed collection of objects.



- Resilient: If data in memory is lost, it can be recreated.
 - Distributed: Processed across the cluster.
 - Dataset: Initial data can come from a source such as a file, or it can be created programmatically.
- RDDs are the fundamental unit of data in Spark.
 - Most Spark programming consists of performing operations on RDDs.

Creating RDD

Python

```
lines = sc.parallelize(["workshop", "spark"])
```

Scala

```
val lines = sc.parallelize(List("workshop", "spark"))
```

Java

```
JavaRDD<String> lines = sc.parallelize(Arrays.asList("workshop", "spark"))
```

Python

```
lines = sc.textFile("/path/to/file.txt")
```

Scala

```
val lines = sc.textFile("/path/to/file.txt")
```

Java

```
JavaRDD<String> lines = sc.textFile("/path/to/file.txt")
```

RDD persistence

MEMORY_ONLY

MEMORY_AND_DISK

MEMORY_ONLY_SER

MEMORY_AND_DISK_SER

DISK_ONLY

MEMORY_ONLY_2

MEMORY_AND_DISK_2

OFF_HEAP

Working with RDDs

RDDs

- RDDs can hold any serializable type of element.
 - o Primitive types: integers, characters, Booleans.
 - o Sequence types: string, lists, arrays, tuples, dictionaries.
 - o Scala/Java Objects (if serializable).
 - o Mixed types.
- Some RDDs are specialized and have additional functionality.
 - o Pair RDDs.
 - o RDDs consisting of key-value pairs.
 - o Double RDDs.
 - o RDDs consisting of numeric data.

Creating RDDs from Collections

```
myData = ["Alice","Carlos","Frank","Barbara"]  
> myRdd = sc.parallelize(myData)  
> myRdd.take(2) ['Alice', 'Carlos']
```

Creating RDDs from Text Files

```
-sc.textFile("myfile.txt")  
-sc.textFile("mydata/")  
-sc.textFile("mydata/*.log")  
-sc.textFile("myfile1.txt,myfile2.txt")
```

- Each line in each file is a separate record in the RDD.
- File are referenced by absolute or relative URI:
 - o Absolute URI:
 - file:/home/training/myfile.txt
 - hdfs://nnhost/loudacre/myfile.txt

Example: Multi-RDD Transformations

rdd1	rdd2
Chicago	San Francisco
Boston	Boston
Paris	Amsterdam
San Francisco	Mumbai
Tokyo	McMurdo Station

`rdd1.subtract(rdd2)`

Tokyo
Paris
Chicago

`rdd1.zip(rdd2)`

(Chicago, San Francisco)
(Boston, Boston)
(Paris, Amsterdam)
(San Francisco, Mumbai)
(Tokyo, McMurdo Station)

rdd1	rdd2
Chicago	San Francisco
Boston	Boston
Paris	Amsterdam
San Francisco	Mumbai
Tokyo	McMurdo Station

`rdd1.union(rdd2)`

Chicago
Boston
Paris
San Francisco
Tokyo
San Francisco
Boston
Amsterdam
Mumbai
McMurdo Station

`rdd1.intersection(rdd2)`

Boston
San Francisco

Some Other General RDD Operations

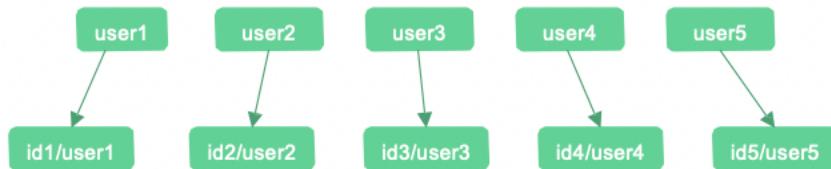
- Other RDD operations:
 - o `first` returns the first element of the RDD.
 - o `foreach` applies a function to each element in an RDD.
 - o `top(n)` returns the largest n elements using natural ordering.
- Sampling operations:
 - o `sample` creates a new RDD with a sampling of elements.
 - o `take` Sample returns an array of sampled elements.

Other data structure in Spark

- Paired RDD.
- Data Frame.
- Data Set.

Paired RDD

- Paired RDD = an RDD of key/value pairs



Pair RDDs

- Pair RDDs are a special form of RDD:
 - o Each element must be a key-value pair (a two-element tuple).
 - o Keys and values can be any type.
- Why?
 - o Use with map-reduce algorithms.
 - o Many additional functions are available for common data processing needs.
 - o Such as sorting, joining, grouping, counting.

Pair RDD
(key1, value1)
(key2, value2)
(key3, value3)
...

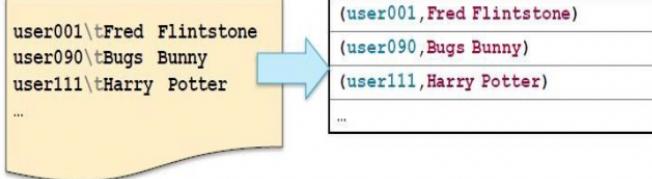
Creating Pair RDDs

- The first step in most workflows is to get the data into key/value form.
 - o What should the RDD be keyed on?
 - o What is the value?
- Commonly used functions to create pair RDDs.
 - o map
 - o flatMap / flatMapValues
 - o keyBy

Example: A simple Pair RDD

Example: Create a pair RDD from a tab-separated file

```
> val users = sc.textFile(file).  
map(line => line.split('\t')).  
map(fields => (fields(0), fields(1)))
```



Example: Keying Web Logs by User ID

```
> sc.textFile(logfile).  
keyBy(line => line.split(' ') (2))
```

User ID

```
56.38.234.188 - 99788 "GET /KBDOC-00157.html HTTP/1.0" ...  
56.38.234.188 - 99788 "GET /theme.css HTTP/1.0" ...  
203.146.17.59 - 25254 "GET /KBDOC-00230.html HTTP/1.0" ...  
...
```

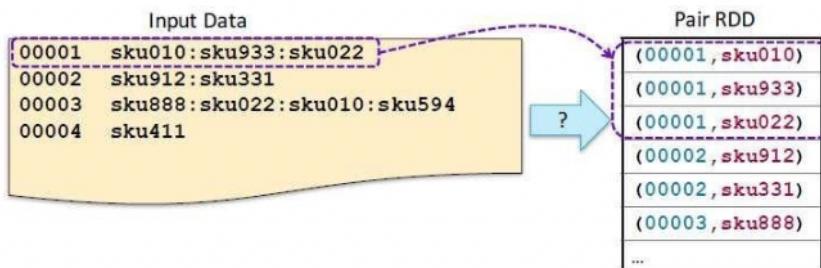
An arrow points from the User ID section to a blue box containing the resulting Pair RDD:

(99788, 56.38.234.188 - 99788 "GET /KBDOC-00157.html...")
(99788, 56.38.234.188 - 99788 "GET /theme.css...")
(25254, 203.146.17.59 - 25254 "GET /KBDOC-00230.html...")
...

Mapping Single Rows to Multiple Pairs

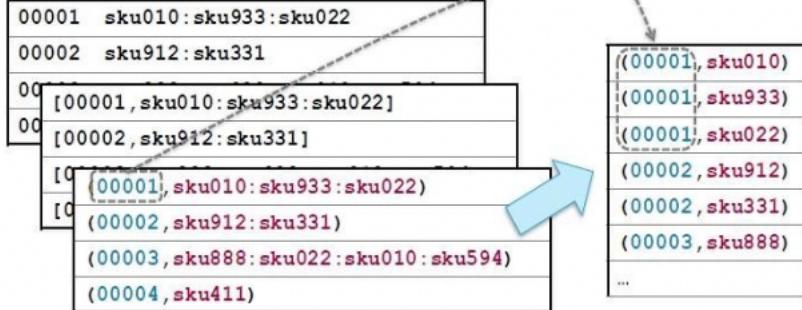
How would you do this?

- Input: order numbers with a list of SKUs in the order
- Output: `order` (key) and `sku` (value)



Answer : Mapping Single Rows to Multiple Pairs

```
> sc.textFile(file) \
    .map(lambda line: line.split('\t')) \
    .map(lambda fields: (fields[0], fields[1])) \
    .flatMapValues(lambda skus: skus.split(':'))
```



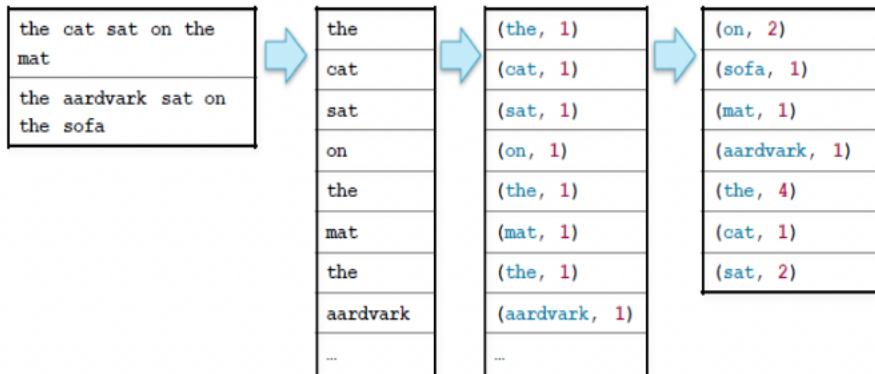
MapReduce

- MapReduce is a common programming model.
 - o Easily applicable to distributed processing of large data sets.
- Hadoop MapReduce is the major implementation.
 - o Somewhat limited.
 - o Each job has one Map phase, one Reduce phase.
 - o Job output is saved to files.
- Spark implements map-reduce with much greater flexibility.
 - o Map and reduce functions can be interspersed.
 - o Results can be stored in memory.
 - o Operations can easily be chained.

MapReduce in Spark

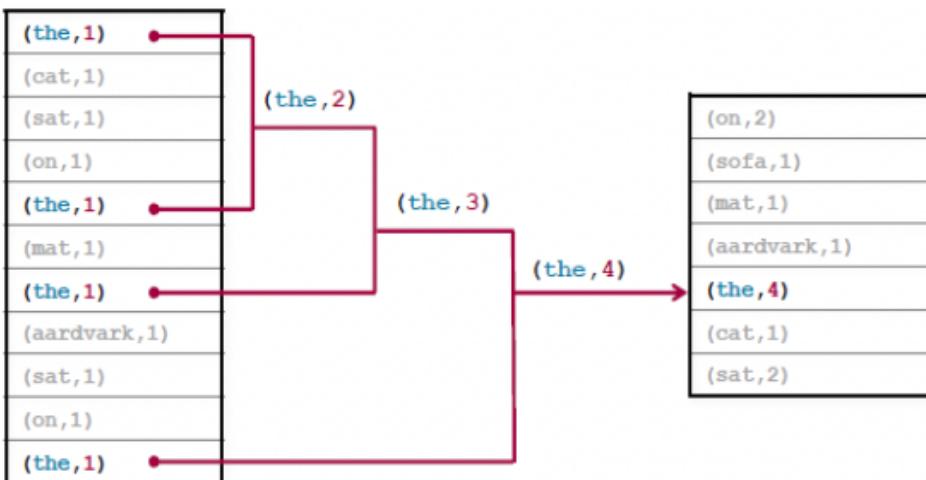
- MapReduce in Spark works on pair RDDs.
- Map phase:
 - o Operates on one record at a time.
 - o “Maps” each record to zero or more new records.
 - o Examples: `map`, `flatMap`, `filter`, `keyBy`.
- Reduce phase:
 - o Works on map output.
 - o Consolidates multiple records.
 - o Examples: `reduceByKey`, `sortByKey`, `mean`.

Example: WordCount



reduceByKey

- The function passed to reduceByKey combines values from two keys.
- Function must be binary.



```
val counts = sc.textFile(file).flatMap(line => line.split(' '))
    .map(word => (word, 1))
    .reduceByKey((v1, v2) => v1 + v2)
```

Or:

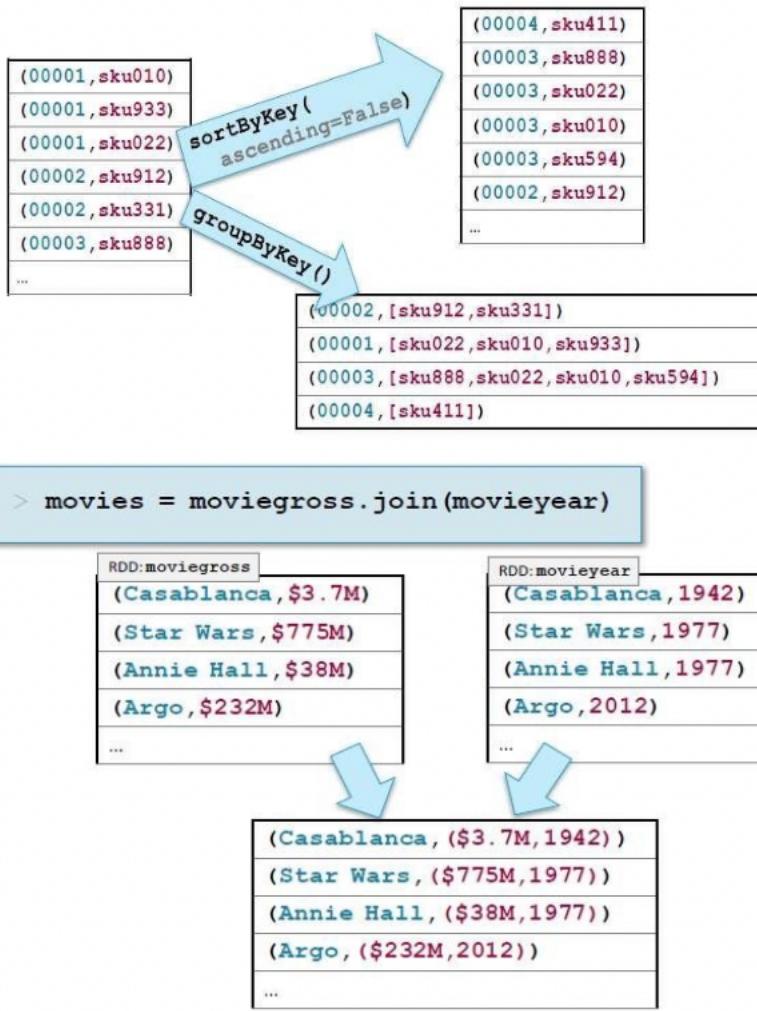
```
val counts = sc.textFile(file).flatMap(_.split('\n')).map(_,_).reduceByKey(_ + _)
```

Pair RDD Operations

- In addition to map and reduceByKey operations, Spark has several operations specific to pair RDDs.
- Examples:
 - o `countByKey` returns a map with the count of occurrences of each key.

- `groupByKey` groups all the values for each key in an RDD.
- `sortByKey` sorts in ascending or descending order.
- `join` returns an RDD containing all pairs with matching keys from two RDDs.

Example: Pair RDD Operations



Other Pair Operations

- Some other pair operations:
 - `keys` return an RDD of just the keys, without the values.
 - `values` return an RDD of just the values, without the keys.
 - `lookup(key)` returns the value(s) for a key.
 - `leftOuterJoin, rightOuterJoin, fullOuterJoin` join two RDDs, including keys defined in the left, right or RDD respectively.
 - `mapValues, flatMapValues` execute a function on just the values, keeping the key the same.

Data Frames and Apache Spark SQL

What is Spark SQL?

- What is Spark SQL?
 - o Spark module for structured data processing.
 - o Replaces Spark (a prior Spark module, now deprecated).
 - o Built on top of core Spark.
- What does Spark SQL provide?
 - o The DataFrame API – a library for working with data as tables.
 - o Defines DataFrame containing rows and columns.
 - o DataFrames are the focus of this chapter.
 - o Catalyst Optimizer – an extensible optimization framework.
 - o A SQL engine and command line interface.

SQL context

- The main Spark SQL entry point is a SQL context object.
 - o Requires a `SparkContext` object.
 - o The SQL context in Spark SQL is similar to Spark context in core Spark.
- There are two implementations:
 - o `SQLContext`.
 - o Basic implementation.
 - o `HiveContext`.
 - o Reads and writes Hive/HCatalog tables directly.
 - o Supports full HiveQL language.
 - o Requires the Spark application be linked with Hive libraries.
 - o Cloudera recommends using `HiveContext`.

Creating a SQL Context

- The Spark shell creates a `HiveContext` instance automatically.
 - o Call `sqlContext`.
 - o You will need to create one when writing a Spark application.
 - o Having multiple SQL context objects is allowed.
- A SQL context object is created based on the Spark context.

Language: Scala

```
import org.apache.spark.sql.hive.HiveContext
val sqlContext = new HiveContext(sc)
import sqlContext.implicits._
```

DataFrames

- DataFrames are the main abstraction in SparkSQL.

- Analogous to RDDs in core Spark.
- A distributed collection of structured data organized into Named columns.
- Built on a base RDD containing **Row** objects.

Creating a DataFrame from a Data Source

- `sqlContext.read` returns a `DataFrameReader` object.
- `DataFrameReader` provides functionality to load data into a `DataFrame`.
- Convenience functions:

-json(filename)

-parquet(filename)

-orc(filename)

-table(hive-tablename)

-jdbc(url,table,options)

Example: Creating a DataFrame from a JSON file



Example: Creating a DataFrame from Hive/Impala Table



Loading from a Data Source Manually

§ You can specify settings for the DataFrameReader

- format**: Specify a data source type
- option**: A key/value setting for the underlying data source
- schema**: Specify a schema instead of inferring from the data source

§ Then call the generic base function **load**

```
sqlContext.read.  
  format("com.databricks.spark.avro") .  
  load("/loudacre/accounts_avro")  
  
sqlContext.read.  
  format("jdbc") .  
  option("url", "jdbc:mysql://localhost/loudacre") .  
  option("dbtable", "accounts") .  
  option("user", "training") .  
  option("password", "training") .  
  load()
```

Data Sources

§ Spark SQL 1.6 built-in data source types

- table**
- json**
- parquet**
- jdbc**
- orc**

§ You can also use third party data source libraries, such as

- Avro (included in CDH)
- HBase
- CSV
- MySQL
- and more being added all the time

DataFrame Basic Operations

- § Basic operations deal with DataFrame metadata (rather than its data)
- § Some examples
 - –schema returns a schema object describing the data
 - –printSchema displays the schema as a visual tree
 - –cache / persist persists the DataFrame to disk or memory
 - –columns returns an array containing the names of the columns
 - –dtypes returns an array of (column name,type) pairs
 - –explain prints debug information about the DataFrame to the console

Language: Scala

```
> val peopleDF = sqlContext.read.json("people.json")
> peopleDF.dtypes.foreach(println)
(age, LongType)
(name, StringType)
(pcode, StringType)
```

DataFrame Actions

§ Some DataFrame actions

- collect returns all rows as an array of Row objects
- take(*n*) returns the first *n* rows as an array of Row objects
- count returns the number of rows
- show(*n*) displays the first *n* rows
(default=20)

Language: Scala

```
> peopleDF.count()
res7: Long = 5
> peopleDF.show(3)
age  name   pcode
null Alice  94304
30  Brayden 94304
19  Carla   10036
```

DataFrame Queries

§ DataFrame query methods return new DataFrames

- Queries can be chained like transformations

§ Some query methods

- **distinct** returns a new DataFrame with distinct elements of this DF
- **join** joins this DataFrame with a second DataFrame
 - Variants for inside, outside, left, and right joins
- **limit** returns a new DataFrame with the first **n** rows of this DF
- **select** returns a new DataFrame with data from one or more columns of the base DataFrame
- **where** returns a new DataFrame with rows meeting specified query criteria (alias for **filter**)

DataFrame Query Strings

- Some query operations take strings containing simple query expressions

- Such as **select** and **where**

- Example: **select**

age	name	pcode
null	Alice	94304
30	Brayden	94304
19	Carla	10036
46	Diana	null
null	Étienne	94104

peopleDF.
select("age")

age
null
30
19
46
null

peopleDF.
select("name", "age")

name	age
Alice	null
Brayden	30
Carla	19
Diana	46
Étienne	null

Querying DataFrames using Columns

§ Columns can be referenced in multiple ways

▪ Scala

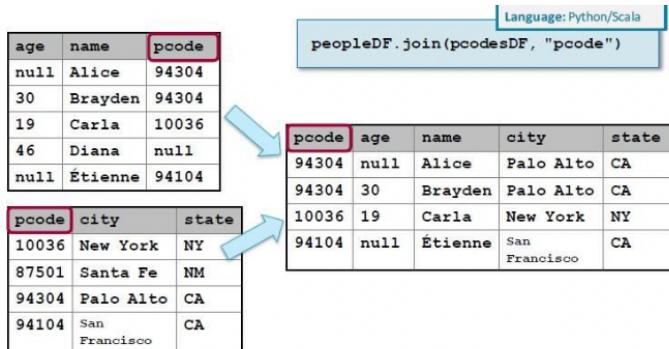
```
val ageDF = peopleDF.select(peopleDF("age"))  
  
val ageDF = peopleDF.select($"age")
```

age	name	pcode
null	Alice	94304
30	Brayden	94304
19	Carla	10036
46	Diana	null
null	Étienne	94104

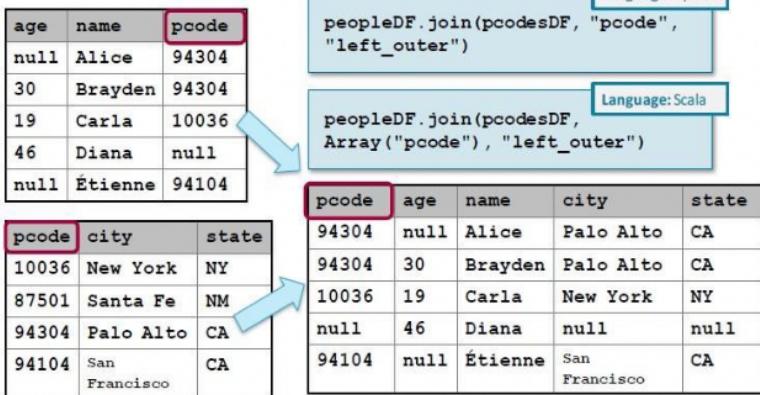
age
null
30
19
46
null

Joining DataFrames

§ A basic inner join when join column is in both DataFrames



▪ Specify type of join as inner (default), outer, left_outer, right_outer, or leftsemi

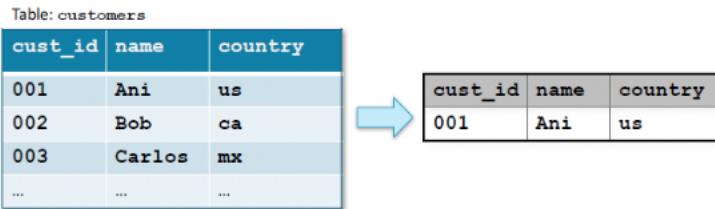


SQL Queries

§ When using **HiveContext**, you can query Hive/Impala tables using **HiveQL**

- Returns a DataFrame

```
Language: Python/Scala  
sqlContext.  
    sql("""SELECT * FROM customers WHERE name LIKE "A%" """)
```



Saving DataFrames

- § Data in DataFrames can be saved to a data source
- § Use **DataFrame.write** to create a **DataStreamWriter**
- § **DataStreamWriter** provides convenience functions to externally save the data represented by a DataFrame
 - jdbc** inserts into a new or existing table in a database
 - json** saves as a JSON file
 - parquet** saves as a Parquet file
 - orc** saves as an ORC file
 - text** saves as a text file (string data in a single column only)
 - saveAsTable** saves as a Hive/Impala table (**HiveContext** only)

```
Language: Python/Scala  
peopleDF.write.saveAsTable("people")
```

Options for Saving DataFrames

- § **DataStreamWriter** option methods
 - format** specifies a data source type
 - mode** determines the behavior if file or table already exists:
 - overwrite, append, ignore or error (default is error)
 - partitionBy** stores data in partitioned directories in the form *column=value* (as with Hive/Impala partitioning)
 - options** specifies properties for the target data source
 - save** is the generic base function to write the data

```
Language: Python/Scala  
peopleDF.write.  
    format("parquet") .  
    mode("append") .  
    partitionBy("age") .  
    saveAsTable("people")
```

DataFrames and RDDs

§ DataFrames are built on RDDs

- Base RDDs contain **Row** objects
- Use **rdd** to get the underlying RDD

```
peopleRDD = peopleDF.rdd
```

peopleDF		
age	name	pcode
null	Alice	94304
30	Brayden	94304
19	Carla	10036
46	Diana	null
null	Étienne	94104

peopleRDD		
Row[null, Alice, 94304]		
Row[30, Brayden, 94304]		
Row[19, Carla, 10036]		
Row[46, Diana, null]		
Row[null, Étienne, 94104]		

§ Row RDDs have all the standard Spark actions and transformations

- Actions: **collect**, **take**, **count**, and so on
- Transformations: **map**, **flatMap**, **filter**, and so on

§ Row RDDs can be transformed into pair RDDs to use map-reduce methods

§ DataFrames also provide convenience methods (such as **map**, **flatMap**, and **foreach**) for converting to RDDs

Working with Row Objects

- Use **Array**-like syntax to return values with type **Any**
- **row(n)** returns element in the *n*th column
- **row.fieldIndex("age")** returns index of the **age** column
- Use methods to get correctly typed values
- **row.getAs[Long]("age")**
- Use type-specific **get** methods to return typed values
- **row.getString(n)** returns *n*th column as a string
- **row.getInt(n)** returns *n*th column as an integer
- And so on

Example: Extracting Data from Row Objects

Extract data from Row objects

```
peopleRDD = peopleDF \
    .map(lambda row: (row.pcode, row.name))
peopleByPCode = peopleRDD \
    .groupByKey()
```

Language: Python

```
val peopleRDD = peopleDF.
    map(row =>
        (row(row.fieldIndex("pcode")),
         row(row.fieldIndex("name"))))
val peopleByPCode = peopleRDD.
    groupByKey()
```

Language: Scala

Row[null, Alice, 94304]
Row[30, Brayden, 94304]
Row[19, Carla, 10036]
Row[46, Diana, null]
Row[null, Étienne, 94104]

(94304, Alice)
(94304, Brayden)
(10036, Carla)
(null, Diana)
(94104, Étienne)

(null, [Diana])
(94304, [Alice, Brayden])
(10036, [Carla])
(94104, [Étienne])

Converting RDDs to DataFrames

You can also create a DF from an RDD using `createDataFrame`

```
import org.apache.spark.sql.types._
import org.apache.spark.sql.Row
val schema = StructType(Array(
    StructField("age", IntegerType, true),
    StructField("name", StringType, true),
    StructField("pcode", StringType, true)))
val rowrdd = sc.parallelize(Array(Row(40, "Abram", "01601"),
                                 Row(16, "Lucia", "87501")))
val mydf = sqlContext.createDataFrame(rowrdd, schema)
```

Language: Scala

Working with SparkRDDs, Pair-RDDs

RDD Operations

Transformations	Actions
map()	count()
flatMap()	collect()
filter()	first(), top(n)
union()	take(n), takeOrdered(n)
intersection()	countByValue()
distinct()	reduce()
groupByKey()	foreach()
reduceByKey()	...
sortByKey()	
join()	
...	

Lambda Expression

PySpark WordCount example:

```
input_file = sc.textFile("/path/to/text/file")
map = input_file.flatMap(lambda line: line.split(" "))
    .map(lambda word: (word, 1))
counts = map.reduceByKey(lambda a, b: a + b)
counts.saveAsTextFile("/path/to/output/")
```

lambda arguments: expression

PySpark RDD API

<https://spark.apache.org/docs/latest/api/python/pyspark.html#pyspark.RDD>

`map(f, preservesPartitioning=False)`

[source]

Return a new RDD by applying a function to each element of this RDD.

```
>>> rdd = sc.parallelize(["b", "a", "c"])
>>> sorted(rdd.map(lambda x: (x, 1)).collect())
[('a', 1), ('b', 1), ('c', 1)]
```

`flatMap(f, preservesPartitioning=False)`

[source]

Return a new RDD by first applying a function to all elements of this RDD, and then flattening the results.

```
>>> rdd = sc.parallelize([2, 3, 4])
>>> sorted(rdd.flatMap(lambda x: range(1, x)).collect())
[1, 1, 1, 2, 2, 3]
>>> sorted(rdd.flatMap(lambda x: [(x, x), (x, x)]).collect())
[(2, 2), (2, 2), (3, 3), (3, 3), (4, 4), (4, 4)]
```

Practice with flight data

Data: airports.dat (<https://openflights.org/data.html>)

[*Airport ID, Name, City, Country, IATA, ICAO, Latitude, Longitude, Altitude, Timezone, DST, Tz database, Type, Source*]

Try to do somethings:

- Create RDD from textfile
- Count the number of airports
- Filter by country
- Group by country
- Count the number of airports in each country

• **Data: airports.dat (<https://openflights.org/data.html>)**

[*Airport ID, Name, City, Country, IATA, ICAO, Latitude, Longitude, Altitude, Timezone, DST, Tz database, Type, Source*]

• **Data: routes.dat**

[*Airline, Airline ID, Source airport, Source airport ID, Destination airport, Destination airport ID, Codeshare, Stops, Equipment*]

Try to do somethings:

- Join 2 RDD
- Count the number of flights arriving in each country

Working with DataFrame and Spark SQL

Creating a DataFrame

```
%pyspark
from pyspark.sql import *

Employee = Row("firstName", "lastName", "email", "salary")

employee1 = Employee('Basher', 'armbrust', 'bash@edureka.co', 100000)
employee2 = Employee('Daniel', 'meng', 'daniel@stanford.edu', 120000 )
employee3 = Employee('Muriel', None, 'muriel@waterloo.edu', 140000 )
employee4 = Employee('Rachel', 'wendell', 'rach_3@edureka.co', 160000 )
employee5 = Employee('Zach', 'galifianakis', 'zach_g@edureka.co', 160000 )

employees = [employee1,employee2,employee3,employee4,employee5]

print(Employee[0])
print(employees)

df = spark.createDataFrame(employees)
df.show()
```

From CSV file:

```
%pyspark
flightData2015 = spark\
    .read\
    .option("inferSchema", "true")\
    .option("header", "true")\
    .csv("/usr/zeppelin/module9/2015-summary.csv")

flightData2015.show()
```

From RDD:

```
%pyspark
from pyspark.sql import *
list = [('Ankit',25),('Jalfaizy',22),('saurabh',20),('Bala',26)]
rdd = sc.parallelize(list)
people = rdd.map(lambda x: Row(name=x[0], age=int(x[1])))
df = spark.createDataFrame(people)

df.show()
```

DataFrame APIs

- **DataFrame**: show(), collect(), createOrReplaceTempView(), distinct(), filter(), select(), count(), groupBy(), join()...
- **Column**: like()
- **Row**: row.key, row[key]
- **GroupedData**: count(), max(), min(), sum(), ...

<https://spark.apache.org/docs/latest/api/python/pyspark.sql.html>

SparkSQL

- Create a temporary view
- Query using SQL syntax

```
%pyspark
flightData2015.createOrReplaceTempView("flight_data_2015")

maxSql = spark.sql("""
SELECT DEST_COUNTRY_NAME, sum(count) as destination_total
FROM flight_data_2015
GROUP BY DEST_COUNTRY_NAME
ORDER BY sum(count) DESC
LIMIT 5
""")

maxSql.show()
```

Table of Contents

<i>Spark Architecture</i>	1
<i>Easy Ways to Run Spark</i>	1
<i>Supported Languages</i>	1
<i>RDD</i>	1
<i>Creating RDD</i>	2
<i>RDD persistence</i>	2
<i>Working with RDDs</i>	3
<i>RDDs</i>	3
<i>Creating RDDs from Collections</i>	3
<i>Creating RDDs from Text Files</i>	3
<i>Example: Multi-RDD Transformations</i>	4
<i>Some Other General RDD Operations</i>	5
<i>Other data structure in Spark</i>	5
<i>Paired RDD</i>	5
<i>Pair RDDs</i>	5
<i>Creating Pair RDDs</i>	5
<i>Example: A simple Pair RDD</i>	6
<i>MapReduce</i>	7
<i>MapReduce in Spark</i>	7
<i>Example: WordCount</i>	8
<i>reduceByKey</i>	8
<i>Pair RDD Operations</i>	8
<i>Example: Pair RDD Operations</i>	9
<i>Other Pair Operations</i>	9
<i>Data Frames and Apache Spark SQL</i>	10
<i>What is Spark SQL?</i>	10
<i>SQL context</i>	10
<i>Creating a SQL Context</i>	10
<i>DataFrames</i>	10
<i>Creating a DataFrame from a Data Source</i>	11

<i>Example: Creating a DataFrame from a JSON file</i>	11
<i>Example: Creating a DataFrame from Hive/Impala Table</i>	11
<i>Loading from a Data Source Manually</i>	12
<i>Data Sources</i>	12
<i>DataFrame Basic Operations</i>	13
<i>DataFrame Actions</i>	13
<i>DataFrame Queries</i>	14
<i>DataFrame Query Strings</i>	14
<i>Querying DataFrames using Columns</i>	15
<i>Joining DataFrames</i>	15
<i>SQL Queries</i>	16
<i>Saving DataFrames</i>	16
<i>Options for Saving DataFrames</i>	16
<i>DataFrames and RDDs</i>	17
<i>Working with Row Objects</i>	17
<i>Example: Extracting Data from Row Objects</i>	18
<i>Converting RDDs to DataFrames</i>	18
<i>Working with SparkRDDs, Pair-RDDs</i>	19
<i>RDD Operations</i>	19
<i>Lambda Expression</i>	19
<i>PySpark RDD API</i>	19
<i>Practice with flight data</i>	20
<i>Working with DataFrame and Spark SQL</i>	20
<i>Creating a DataFrame</i>	20
<i>DataFrame APIs</i>	21
<i>SparkSQL</i>	21