

It is a group or much more than a programming paradigm but an variety of packages, libraries, and systems built on top of the Co

Spark Core consists of two APIs. The Unstructured and Structured Unstructured API is Spark's lower level set of APIs including RDDs, Distributed Datasets (RDDs), Accumulators, and Broadcast variables. The Structured API consists of DataFrames, Datasets, Spark SQL and MLlib. The DataFrame API is an interface that most users should use. The difference between the two APIs is that the DataFrame API is optimized to work with structured data in a spreadsheet-like interface while the other is meant for manipulation of raw java objects. Outside of the core Spark library, there are many tools and modules. There is MLLib for performing machine learning, the GraphX module for performing graph processing, and SparkR for working with Spark clusters from the R language.

We will cover all of these tools in due time however this chapter will focus on the core Spark library.

processing. Single machines simply cannot have enough power to perform computations on huge amounts of information (or they don't have time to wait for the computation to finish). A *cluster*, or group of machines, pools the resources of many machines together. Now having many machines alone is not powerful, you need a framework to coordinate them and manage their resources across them. Spark is a tool for just that, managing and coordinating many machines.

# Spark Applications

Spark Applications consist of a *driver* process and a set of *executors*. The driver process, Figure 1-2, sits on the *driver node* and is responsible for three things: maintaining information about the application, responding to a user's program, and analyzing, distributing work across the executors. As suggested by figure 1, the driver process is absolutely essential - it's the heart of a Spark Application. It maintains all relevant information during the lifetime of the application.



An executor is responsible for two things: executing code assigned by the driver and reporting the state of the computation back to the driver.

The last piece relevant piece for us is the cluster manager. The cluster manager controls physical machines and allocates resources to Spark applications. There can be one of several core cluster managers: Spark's standalone cluster manager, YARN, or Mesos. This means that there can be multiple applications running on a cluster at the same time.

## NOTE:

Spark, in addition to its cluster mode, also has a *local mode*. Recall that the driver and executors are processes? This means that Spark can run where these processes live. In *local mode*, these processes run on an *individual computer instead of a cluster*. See figure 1-3 for a high-level diagram of this architecture. This is the easiest way to get started and what the demonstrations in this book should run on.

## Scala

Spark is primarily written in Scala, making it Spark's "default" language. This book will include examples of Scala wherever there are code snippets.

## Python

Python supports nearly everything that Scala supports. This book will include numerous Python API examples wherever possible.

Even though Spark is written in Scala, Spark's authors have been careful to ensure that you can write Spark code in Java. This book will focus on Scala but will provide Java examples where relevant.

## SQL

Spark supports user code written in ANSI 2003 Compliant SQL. This makes it easy for analysts and non-programmers to leverage the big data features of Spark. This book will include numerous SQL examples.

## R

Spark supports the execution of R code through a project called

will cover this in the Ecosystem section of the book along with interesting projects that aim to do the same thing like [Sparklyr](#).

architecture because at this point it's not necessary to get us closer to our own Spark code. The key points are that:

- Spark has some cluster manager that maintains an understanding of resources available.
- The driver process is responsible for executing our driver commands across the executors in order to complete our tasks.
- There are two modes that you can use, cluster mode (on multiple machines) and local mode (on a single machine).

Now in the previous chapter we talked about what you need to do to work with Spark by setting your Java, Scala, and Python versions. Now to start Spark's local mode, this means running `./bin/spark-shell`. Once you start that you will see a console, into which you can enter commands. If you would like to work in Python you would run `./bin/pyspark`.

From the beginning of this chapter we know that we leverage a driver process to maintain our Spark Application. This driver process manifests the user as something called the `SparkSession`. The `SparkSession` is the entrance point to executing code in Spark, in any language, and is the facing part of a Spark Application. In Scala and Python the variable is available as `spark` when you start up the Spark console. Let's go ahead and look at the `SparkSession` in both Scala and/or Python.

```
%scala
```

```
spark
```

```
%python
```

```
spark
```

In Scala, you should see something like:

```
res0: org.apache.spark.sql.SparkSession = org.apache.spark.sql.SparkSession@7efda4c1cccd
```

In Python you'll see something like:

```
<pyspark.sql.session.SparkSession at 0x7efda4c1cccd>
```

Now you need to understand how to submit commands to the Spark Session. Let's do that by performing one of the simplest tasks that we can.

range of numbers. This range of numbers is just like a named column in a spreadsheet.

```
%scala  
val myRange = spark.range(1000).toDF("number")  
  
%python  
myRange = spark.range(1000).toDF("number")
```

You just ran your first Spark code! We created a DataFrame with 1000 rows containing values from 0 to 999. This range of numbers represents a *distributed collection*. Running on a cluster, each partition of this range of numbers would exist on a different executor. You'll notice that the value of *myRange* is a DataFrame, let's introduce DataFrames!

## DataFrames

A *DataFrame* is a table of data with rows and columns. We call columns and their types a *schema*. A simple analogy would be a table with named columns. The fundamental difference is that while a table sits on one computer in one specific location, a Spark DataFrame potentially thousands of computers. The reason for putting the data on more than one computer should be intuitive: either the data is too large to fit on one machine or it would simply take too long to perform that computation on one machine.

The DataFrame concept is not unique to Spark. The R Language has a similar concept as do certain libraries in the Python programming language. Python/R DataFrames (with some exceptions) exist on one machine or multiple machines. This limits what you can do with a given DataFrame since python and R are limited to the resources that exist on that specific machine. Since Spark has language interfaces for both Python and R, it's easy to convert between them. You can convert Pandas (Python) DataFrames to Spark DataFrames and vice versa. You can also convert Spark DataFrames to Pandas DataFrames (in R).

### Note

Spark has several core abstractions: Datasets, DataFrames, SQL Databases, and Resilient Distributed Datasets (RDDs). These abstractions all represent distributed collections of data however they have different interfaces for working with that data. The easiest and most efficient are DataFrames and RDDs. These are available in all languages. We cover Datasets in Section II, DataFrames in depth in Section III Chapter 2 and 3. The following covers all of the core abstractions.

## Partitions

In order to leverage the resources of the machines in cluster, we break up the data into chunks, called partitions. A *partition* is a collection of data that sit on one physical machine in our cluster. A DataFrame can have more partitions.

Spreadsheet on a single machine

When we perform some computation, Spark will operate on each partition parallel unless an operation calls for a *shuffle*, where multiple partitions need to share data. Think about it this way, if you need to run some errands, you typically have to do those one by one, or serially. What if you could give one errand to a worker who would then complete that task and return it back to you? In that scenario, the key is to break up errands efficiently so you can get as much work done in as little time as possible. In the context of Spark, an “errand” is equivalent to computation + data and a “worker” is equivalent to an executor.

Now with DataFrames, we do not manipulate partitions individually, instead we give us the `DataFrame` interface for doing that. Now when we write code, you’ll notice there was no list of numbers, only a type signature because Spark organizes computation into two categories, `transformations` and `actions`. When we create a DataFrame, we perform a transformation.

## Transformations

In Spark, the core data structures are *immutable* meaning they cannot be changed once created. This might seem like a strange concept at first. If you cannot change it, how are you supposed to use it? In order to “change” a DataFrame you will have to instruct Spark how you would like the DataFrame you have into the one that you want. These instructions are called *transformations*. Transformations are how you, as user, specify what you would like to transform the DataFrame you currently have to the DataFrame that you want to have. Let’s show an example. To computer whether a number is divisible by two, we use the modulo operation to see if there is anything left over from dividing one number by another.

results that do not result in zero. We can specify this filter using the `where` clause.

```
%scala  
val divisBy2 = myRange.where("number % 2 = 0")
```

```
%python  
divisBy2 = myRange.where("number % 2 = 0")
```

**Note**

DataFrame into a table.

These operations create a new DataFrame but do not execute any code. The reason for this is that DataFrame transformations do not trigger any execution. Instead, they are lazily evaluated.

## Lazy Evaluation

*Lazy evaluation* means that Spark will wait until the very last moment to execute your transformations. In Spark, instead of modifying the DataFrame directly, we build up a plan of the transformations that we would like to perform. By waiting for the last minute to execute your code, can try and run as efficiently as possible across the cluster.

## ACTIONS

To trigger the computation, we run an *action*. An *action* instructs compute a result from a series of *transformations*. The simplest is `count` which gives us the total number of records in the DataFrame.

```
%scala  
divisBy2.count()  
  
%python  
divisBy2.count()
```

We now see a result! There are 500 numbers divisible by two from (big surprise!). Now `count` is not the only action. There are three actions: actions to view data in the console, actions to collect data objects in the respective language, and actions to write to output.

## SPARK UI

During Spark's execution of the previous code block, users can monitor progress of their job through the Spark UI. The Spark UI is available at port 4040 of the driver node. If you are running in local mode this would be <http://localhost:4040>. The Spark UI maintains information about our Spark jobs, environment, and cluster state. It's very useful, for tuning and debugging. In this case, we can see one Spark job with one and nine tasks were executed.

In this chapter we will avoid the details of Spark jobs and the point you should understand that a Spark job represents a set of triggered by an individual action. We talk in depth about the breakdown of a Spark job in Section IV.

we using some [flight data](#) from the United States Bureau of Transport Statistics.

We touched briefly on the `SparkSession` as the interface the entire system uses for performing work on the Spark cluster. the `SparkSession` can do simply parallelize an array it can create `DataFrames` directly from `JSON` files. In this case, we will create our `DataFrames` from a `JSON` file that contains some summary flight information collected by the United States Bureau of Transport Statistics. In

What we'll do is start with one specific year and then work up to more years. Let's go ahead and create a `DataFrame` from 2015. To do this, we use the `DataFrameReader` (via `spark.read`) interface, specifying the path.

```
%scala  
  
val flightData2015 = spark  
  .read  
  .json("/mnt/defg/chapter-1-data/json/2015-summary.json")  
  
%python  
  
flightData2015 = spark\
```

You'll see that our two DataFrames (in Scala and Python) each have columns with an unspecified number of rows. Let's take a peek at a new action, `take`, which allows us to view the first couple of rows of a DataFrame. Figure 1-7 illustrates the conceptual actions that we will process. We lazily create the DataFrame then call an action to take values.

```
%scala  
flightData2015.take(2)  
  
%python  
flightData2015.take(2)
```

Remember how we talked about Spark building up a plan, this is a conceptual tool, this is actually what happens under the hood. We can see the actual plan built by Spark by running the `explain` method.

```
flightData2015.explain()  
  
%python  
flightData2015.explain()
```

This just describes reading data from a certain location however as you will start to notice patterns in the explain plans. Without going into much detail at this point, the explain plan represents the logical

transformations Spark will run on the cluster. We can use this to make our code is as optimized as possible. We will not cover that in this chapter, but we will touch on it in the optimization chapter.

#### Note

Remember, we cannot modify this DataFrame by specifying the transform transformation, we can only create a new DataFrame by transforming the previous DataFrame. We can see that even though we're seeing the computation to be completed Spark doesn't yet execute this command, it's just building up a plan. The illustration in figure 1-8 represents what we see in the explain plan for that DataFrame.

Execution Plan

```
%scala  
  
val sortedFlightData2015 = flightData2015.sort("count")  
sortedFlightData2015.explain()  
  
%python  
  
sortedFlightData2015 = flightData2015.sort("count")  
sortedFlightData2015.explain()
```

Now, just like we did before, we can specify an action in order plan.

```
%scala  
  
sortedFlightData2015.take(2)  
  
%python  
  
sortedFlightData2015.take(2)
```

The conceptual plan that we executed previously is illustrated in the following diagram:

Now this planning process is essentially defining *lineage* for the DataFrame, so that at any given point in time Spark knows how to recompute a DataFrame from its dependencies. Given a DataFrame all the way back to a robust data source be it a database or a file. Now that we performed this action, remember that we can go to the Spark UI (port 4040) and see the information about this job and its tasks.

Now hopefully you have grasped the basics but let's just reinforce

## THE FIVE STAGES OF A CLOUD SYSTEM

Let's go ahead and define our DataFrame just like we did before. This time we're going to specify an option for our DataFrameReader.

allow you to control how you read in a given file format and take advantage of some of the structures or information available in case we're going to use two popular options `inferSchema` and

```
%scala  
  
val flightData2015 = spark.read  
  .option("inferSchema", "true")  
  .option("header", "true")  
  .csv("/mnt/defg/chapter-1-data/csv/2015-summary.  
flightData2015
```

```
%python  
  
flightData2015 = spark.read\  
  .option("inferSchema", "true")\  
  .option("header", "true")\  
  .csv("/mnt/defg/chapter-1-data/csv/2015-summary.  
flightData2015
```

After running the code you should notice that we've basically ended up with the same DataFrame that we had when we read in our data from json. The schema is correct looking column names and types. However, we had to be careful when it came to reading in the CSV file as opposed to json because JSON provides a bit more structure than CSVs because JSON has a notion of type.

Looking at them, the `header` option should feel like it makes sense. If the first row in our csv file is the header (column names) and because CSV files are not guaranteed to have this information we must specify it manually. The `inferSchema` option might feel a bit more unfamiliar. JSON obviously provides a bit more structure than CSVs because JSON has a notion of type. To get around this past this by inferring the schema of the csv file we are reading in. Spark can't do this magically, it must scan (read in) some of the data in order to infer the schema, but this saves us from having to specify the types for each column. It also reduces the risk of Spark potentially making an erroneous guess as to what type a column should be.

A discerning reader might notice that the schema returned by our CSV reader does not exactly match that of the json reader.

```
.option("header", "true")
.load("/mnt/defg/chapter-1-data/csv/2015-summary"
.schema

val jsonSchema = spark
.read.format("json")
.load("/mnt/defg/chapter-1-data/json/2015-summar
.schema

println(csvSchema)
println(jsonSchema)

%python
```

```
df1 = pd.read_csv('2015-summary.csv')
```

```
print(csvSchema)
print(jsonSchema)
```

The csv schema:

```
StructType(StructField(DEST_COUNTRY_NAME,StringType,
StructField(ORIGIN_COUNTRY_NAME,StringType,true),
StructField(count,IntegerType,true)))
```

The JSON schema:

```
StructType(StructField(DEST_COUNTRY_NAME,StringType,
StructField(ORIGIN_COUNTRY_NAME,StringType,true),
StructField(count,LongType,true)))
```

For our purposes the difference between a `LongType` and an `IntegerType` is of little consequence however this may be of greater significance in certain scenarios. Naturally we can always explicitly set a schema (rather than inferring it) when we read in data as well. These are just a few options we have when we read in data, to learn more about these options see the chapter on reading and writing data.

```
%scala  
  
val flightData2015 = spark.read  
  .schema(jsonSchema)  
  .option("header", "true")  
  .csv("/mnt/defg/chapter-1-data/csv/2015-summary.  
  
%python  
  
flightData2015 = spark.read\  
  .schema(jsonSchema)\\  
  .option("header", "true")\\  
  .csv("/mnt/defg/chapter-1-data/csv/2015-summary.
```

## DataFrames and SQL

Spark provides another way to query and operate on our DataFrame - that's with SQL! Spark SQL allows you as a user to register any DataFrame or view (a temporary table) and query it using pure SQL. There is no performance difference between writing SQL queries or writing DataFrame code, they both "compile" to the same underlying plan that we see in DataFrame code.

Any DataFrame can be made into a table or view with one simple command:

```
%scala  
flightData2015.createOrReplaceTempView("flight_data")  
  
%python  
flightData2015.createOrReplaceTempView("flight_data")
```

Now we can query our data in SQL. To execute a SQL query, we use the `spark.sql` function (remember `spark` is our `SparkSession` variable). This function is very convenient, because it returns a new DataFrame. While this may seem a bit counter-intuitive at first logic - that a SQL query against a DataFrame returns another DataFrame - it is actually quite powerful. As a user, you can specify transformations in a manner most convenient to you at any given point in time and no longer have to worry about the internal representation of your data.

two explain plans.

%scala

```
val sqlWay = spark.sql("""  
SELECT DEST_COUNTRY_NAME, count(1)  
FROM flight_data_2015  
GROUP BY DEST_COUNTRY_NAME  
""")  
  
val dataFrameWay = flightData2015  
    .groupBy('DEST_COUNTRY_NAME)  
    .count()
```

```
sqlWay.explain  
dataFrameWay.explain
```

%python

```
sqlWay = spark.sql("""
```

transformation, as we are effectively moving down to one row that looks like.

```
// scala or python
spark.sql("SELECT max(count) from flight_data_2015")  
%scala  
  
import org.apache.spark.sql.functions.max  
  
flightData2015.select(max("count")).take(1)  
  
%python  
  
from pyspark.sql.functions import max  
  
flightData2015.select(max("count")).take(1)
```

Let's move onto something a bit more complicated. What are the

query so we'll take it step by step. We will start with a fairly simple SQL aggregation.

```
%scala

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

maxSql.collect()

%python

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

maxSql.collect()
```

Now let's move to the `DataFrame` syntax that is semantically similar but slightly different in implementation and ordering. But, as we mentioned, the underlying plans for both of them are the same. Let's execute them and see their results as a sanity check.

```
%scala
import org.apache.spark.sql.functions.desc

flightData2015
    .groupBy("DEST_COUNTRY_NAME")
    .sum("count")
    .withColumnRenamed("sum(count)", "destination_to")
    .sort(desc("destination_to"))
    .limit(5)
    .collect()

%python
from pyspark.sql.functions import desc
```

```
.limit(5)\n.collect()
```

Now there are 7 steps that take us all the way back to the source. Illustrated below are the set of steps that we perform in “code”. execution plan (the one visible in explain) will differ from what is shown below because of optimizations in physical execution, however, it is as good of a starting point as any. With Spark, we are always working with a directed acyclic graph of transformations resulting in immutable objects. we can subsequently call an action on to see a result.

The first step is to read in the data. We defined the DataFrame previously as a reminder, Spark does not actually read it in until an action is performed on the DataFrame or one derived from the original DataFrame.

The second step is our grouping, technically when we call “groupBy” we are grouping up with a RelationalGroupedDataset which is a fancy name for a DataFrame that has a grouping specified but needs a user to specify an aggregation function before it can be queried further. We can see this by trying to perform an action on the grouped DataFrame (which will not work). We still haven’t performed any computation yet.

of the data.

Therefore the third step is to specify the aggregation. Let's use the aggregation method. This takes as input a column expression or column name. The result of the sum method call is a new DataFrame that it has a new schema but that it does know the type of each column. It is important to reinforce (again!) that no computation has been performed - this is simply another transformation that we've expressed and Spark is responsible for tracing the type information we have supplied.

The fourth step is a simple renaming, we use the `withColumnRenamed` method that takes two arguments, the original column name and the new column name. Of course, this doesn't perform computation - this is just another transformation!

The fifth step sorts the data such that if we were to take results one by one, they would be sorted.

```
destination_total column.
```

You likely noticed that we had to import a function to do this, the `desc` function. You might also notice that `desc` does not return a string. In general, many DataFrame methods will accept Strings (as column names) or expressions. Columns and expressions are actually the same thing.

The final step is just a limit. This just specifies that we only want 10 results. This is just like a filter except that it filters by position (lazily) on the value. It's safe to say that it basically just specifies a DataFrame size.

The last step is our action! Now we actually begin the process of calculating the results of our DataFrame above and Spark will give us back a DataFrame in the language that we're executing. Now to reinforce all of this, let's explain plan for the above query.

```
flightData2015
    .groupBy("DEST_COUNTRY_NAME")
    .sum("count")
```

```
.withColumnRenamed("sum(count)", "destination_to")
.sort(desc("destination_total"))
.limit(5)
.explain()

== Physical Plan ==
TakeOrderedAndProject(limit=5, orderBy=[destination_to])
+- *HashAggregate(keys=[DEST_COUNTRY_NAME#7323], f...
   +- Exchange hashpartitioning(DEST_COUNTRY_NAME#7323)
      +- *HashAggregate(keys=[DEST_COUNTRY_NAME#7323], f...
         +- InMemoryTableScan [DEST_COUNTRY_NAME#7323]
            +- InMemoryRelation [DEST_COUNTRY_NAME#7323]
               +- *Scan csv [DEST_COUNTRY_NAME#7323]
```

While this explain plan doesn't match our exact "conceptual plan", the pieces are there. You can see the limit statement as well as the sort (first line). You can also see how our aggregation happens in two stages of partial\_sum calls. This is because summing a list of numbers is expensive and Spark can perform the sum, partition by partition. Of course, we can also see the InMemoryRelation and Scan calls for how we read in the DataFrame as well.

You are now equipped with the Spark knowledge to writing your own code. In the next chapter we will explore some of Spark's more advanced features.

## SPARK'S STRUCTURED APIs

For our purposes there is a spectrum of types of data. The two ends of the spectrum are **structured** and **unstructured**. Structured and semi-structured data refer to data that have structure that a computer can understand easily. Unstructured data, like a poem or prose, is much harder to understand. Spark's Structured APIs allow for transformations on structured and semi-structured data.

The Structured APIs specifically refer to operations on **DataFrames** and **DataSets** and were created as a high level interface for manipulating big data. This section will cover all the principles of these APIs. Although distinct in the book, the vast majority of these user operations apply to both *batch* as well as *streaming* computations.

## BOX

Before proceeding, let's review the fundamental concepts that we covered in the previous section. Spark is a distributed programming model where the user specifies *transformations* to build up a directed-acyclic-graph of instructions, and *actions* to begin the process of executing that graph of instructions, as by breaking it down into stages and tasks to execute across partitions. The way we store data on which to perform transformations are `DataFrames` and `Datasets`. To create a new `DataFrame`, you call a transformation. To start computation or convert to native types, you call an action.

In Section I, we talked all about DataFrames. Spark has two not “structured” data structures: **DataFrames** and **Datasets**. We will (nuanced) differences shortly but let’s define what they both represent.

To the user, **DataFrames** and **Datasets** are **(distributed) tables** with **columns**. Each column must have the same number of rows as all other columns (although you can use null to specify the lack of a value). **DataFrames** also have type information that tells the user what exists in each column.

To Spark, **DataFrames** and **Datasets** represent by **immutable, lazily evaluated plans** that specify how to perform a series of transformations to produce the correct output. When we perform an action on a **DataFrame** we are **not** performing the actual transformations that represent that **DataFrame**, we are **representing plans of how to manipulate rows and columns to compute the desired result**. Let’s go over rows and column to more precisely define these concepts.

## Schemas

One core concept that differentiates the Structured APIs from the RDD APIs is the concept of a **schema**. A schema defines the column structure of a **DataFrame**. Users can define schemas manually or users can automatically generate a schema from a data source (often called **schema on read**). Now that we know what defines **DataFrames** and **Datasets** and how they get their structure, let’s see an overview of all of the types.

Spark is effectively a programming language of its own. When you perform operations with Spark, it maintains its own type information throughout the execution process. This allows it to perform a wide variety of optimizations during the execution process. These types correspond to the types that Spark supports in each of Scala, Java, Python, SQL, and R. Even if we use Spark's APIs from Python, the majority of our manipulations will operate on **Spark types**, not Python types. For example, the below code does addition in Scala or Python, **it actually performs addition purely in memory**.

```
%scala  
  
val df = spark.range(500).toDF("number")  
df.select(df.col("number") + 10)  
// org.apache.spark.sql.DataFrame = [(number + 10): bigint]  
  
%python  
  
df = spark.range(500).toDF("number")  
df.select(df["number"] + 10)  
# DataFrame[(number + 10): bigint]
```

This is because, as mentioned, **Spark maintains its own type information stored as a schema**, and through some magic in each languages API, it converts an expression in one language to Spark's representation.

## NOTE

There are two distinct APIs within the Structured API. The first API that goes across languages, more commonly referred to as the **DataFrame** or "untyped API". The second API is the "typed API", that is only available to JVM based languages (Scala, Java). This is a bit of a misnomer because the "untyped API" can work with any type, but it only operates on **Spark types at run time**. The "typed API" allows you to define your own types to represent each record.

dataset with “case classes or Java Beans” and types are checked at compile time. Each record (or row) in the “untyped API” contains Row objects that are available across languages and still have only Spark types, not native types. The “typed API” is covered in the Datasets Chapter at the end of Section II. The majority of examples in this chapter cover the “untyped API” however all of this still applies to the typed API.

Notice how the following code produces a dataset of type long, which is an internal Spark type (bigint).

```
%scala  
spark.range(500)
```

Notice how the following code produces a DataFrame with an integer column of type (bigint).

```
%python
```

## Columns

For now, all you need to understand about columns is that they can be of a simple type like an integer or string, a complex types like an array or a null value. Spark tracks all of this type information to you and helps you in many ways that you can transform columns. Columns are discussed extensively in the next chapter but for the most part you can think about Spark columns as columns in a table.

# Rows

There are two ways of getting data into Spark, through Rows and Row objects. Row objects are the most general way of getting data into, and are available in all languages. Each record in a DataFrame type as we can see when we collect the following DataFrames.

```
%scala
```

```
spark.range(2).toDF().collect()
```

```
%python
```

```
spark.range(2).collect()
```

corresponding language specific types. Caveats and details are reader as well to make it easier to reference.

To work with the correct Scala types:

```
import org.apache.spark.sql.types.  
val b = ByteType()
```

To work with the correct Java types you should use the factory methods in the following package:

```
import org.apache.spark.sql.types.DataTypes;  
ByteType x = DataTypes.ByteType();
```

Python types at time have certain requirements (like the listed requirements for `ByteType` below). To work with the correct Python types:

```
from pyspark.sql.types import *  
b = byteType()
```

**Spark Type**

**Scala Value Type**

**Scala Class**

## ENCODERS

Using Spark from Scala and Java allows you to define your own types and use in place of Rows that consist of the above data types. To do this we can use the Encoder. Encoders are only available in Scala, with case classes and JavaBeans. For some types, like `Long`, Spark already includes an Encoder. For instance we can collect a Dataset of type `Long` and convert it back to Scala types back. We will cover encoders in the Datasets chapter.

```
spark.range(2).collect()
```

structured API query from user code to executed code. As an overview, the steps are:

1. Write DataFrame/Dataset/SQL Code
2. If valid code, Spark converts this to a *Logical Plan*
3. Spark transforms this *Logical Plan* to a *Physical Plan*
4. Spark then executes this *Physical Plan* on the cluster

To execute code, we have to write code. This code is then submitted either through the console or via a submitted job. This code then passes through the Catalyst Optimizer which decides how the code should be run and lays out a plan for doing so, before finally the code is run and the results are returned to the user.

This logical plan only represents a set of abstract transformations. It does not refer to executors or drivers, it's purely to convert the user's sequence of operations into the most optimized version. It does this by converting user's *unresolved logical plan*. This unresolved because while your query may be valid, the tables or columns that it refers to may or may not exist. The logical plan is resolved by the catalog, a repository of all table and DataFrame information. The catalog is used to resolve columns and tables in the analyzer. The analyzer may resolve the logical plan if it finds the required table or column name in the catalog. If it can resolve it, this result is passed through the collection of rules, which attempts to optimize the logical plan by pruning down predicates or selections.

## Physical Planning

After successfully creating an optimized logical plan, Spark begins the physical planning process. The physical plan, often called a Spark SQL execution plan, specifies how the logical plan will execute on the cluster by generating different physical execution strategies and comparing them through a cost model. An example of the cost comparison might be choosing how to perform a given join by looking at the physical attributes of a given table (e.g., how many partitions the table has or how big its partitions are.)

Physical ...

Physical planning results in a series of RDDs and transformations. This is why you may have heard Spark referred to as a compiler, it takes DataFrames, Datasets, and SQL and compiles them into RDD transformations for you.

## Execution

Upon selecting a physical plan, Spark runs all of this code over the lower-level programming interface of Spark covered in Part III. Spark then performs further optimizations by, at runtime, generating native code that can remove whole tasks or stages during execution. Finally, the results are returned to the user.

in depth later in this section.

Definitionally, a DataFrame consists of a series of records (like table), that are of type Row, and a number of columns (like columns in spreadsheet) that represent an computation expression that can produce each individual record in the dataset. The schema defines the name and the type of data in each column. The partitioning of the DataFrame defines the layout of the DataFrame or Dataset's physical distribution across partitions. The partitioning scheme defines how that is broken up, this can be based on values in a certain column or non-deterministically.

Let's define a DataFrame to work with.

```
%scala  
  
val df = spark.read.format("json")  
    .load("/mnt/defg/flight-data/json/2015-summary.json")  
  
%python  
  
df = spark.read.format("json")  
    .load("/mnt/defg/flight-data/json/2015-summary.json")
```

We discussed that a DataFrame will have columns, and we use a command to view all of those. We can run the following command in Scala or Python.

```
df.printSchema()
```

The schema ties the logical pieces together and is the starting point to understand DataFrames.

# Schemas

A schema defines the column names and types of a DataFrame. You can let a data source define the schema (called *schema on read*) or you can define it explicitly ourselves.

## NOTE

Deciding whether or not you need to define a schema prior to reading your data depends on your use case. Often times for ad hoc analysis on text file formats like csv or json. However, this can also lead to precision issues like a long type incorrectly set as an integer when reading in a file. When using Spark for production ETL, it is a good idea to define your schemas manually, especially when working with untyped data sources like csv and json because schema inference can vary depending on the type of data that you read in.

Let's start with a simple file we saw in the previous chapter and see how the structured nature of line-delimited JSON define the structure. The file contains flight data from the United States Bureau of Transportation statistics.

```
%scala  
  
spark.read.format("json")  
  .load("/mnt/defg/flight-data/json/2015-summary.json")  
  .schema
```

Scala will return:

```
org.apache.spark.sql.types.StructType = ...  
StructType(StructField(DEST_COUNTRY_NAME,StringType,  
StructField(ORIGIN_COUNTRY_NAME,StringType,true),  
StructField(count,LongType,true))
```

```
%python
```

```
StructField(ORIGIN_COUNTRY_NAME, StringType, true),  
StructField(count, LongType, true)))
```

A schema is a `StructType` made up of a number of fields, `StructField`. Each field has a name, type, and a boolean flag which specifies whether or not the column can contain missing or `null` values. Schemas can also contain other `StructType` (Spark's complex types). We will see this in the next section when we discuss working with complex types.

If the data types in the data (at runtime), do not match the schema, Spark will throw an error.

```
%scala  
  
import org.apache.spark.sql.types.{StructField, StructType}  
  
val myManualSchema = new StructType(Array(  
    new StructField("DEST_COUNTRY_NAME", StringType, true),  
    new StructField("ORIGIN_COUNTRY_NAME", StringType, true),  
    new StructField("count", LongType, false) // just for fun  
))  
  
val df = spark.read.format("json")  
    .schema(myManualSchema)  
    .load("/mnt/defg/flight-data/json/2015-summary.json")
```

Here's how to do the same in Python.

```
.schema(myManualSchema) \  
.load("/mnt/defg/flight-data/json/2015-summary.json")
```

As discussed in the previous chapter, we cannot simply set type language types because Spark maintains its own type information. We will discuss what schemas define, columns.

## Columns and Expressions

To users, columns in Spark are similar to columns in a spreadsheet, data frame, pandas DataFrame. We can select, manipulate, and re-use data from DataFrames and these operations are represented as *expressions*.

To Spark, columns are logical constructions that simply represent values computed on a per-record basis by means of an *expression*. This means that in order to have a real value for a column, we need to have a row, and in order to have a row we need to have a DataFrame. This means that we can only manipulate an actual column outside of a DataFrame, we can only create a logical column's expressions then perform that expression within a DataFrame.

# Columns

There are a lot of different ways to construct and or refer to columns. The two simplest ways are with the `col` or `column` functions. To use these functions, we pass in a column name.

```
%scala  
  
import org.apache.spark.sql.functions.{col, column}  
  
col("someColumnName")  
column("someColumnName")  
  
%python  
  
from pyspark.sql.functions import col, column  
  
col("someColumnName")  
column("someColumnName")
```

We will stick to using `col` throughout this book. As mentioned, a column may or may not exist in our DataFrame. This is because, as mentioned in the previous chapter, columns are not resolved until we compare their names with those we are maintaining in the *catalog*. Column resolution happens in the *analyzer* phase as discussed in the first part of this section.

NOTE

thing as what we have already, namely creating a column, a performance improvement.

```
%scala  
$"myColumn"  
'myColumn
```

The `$` allows us to designate a string as a special string that should be interpreted as a column expression. The tick mark `'` is a special thing called a symbol, it is a specific construct of referring to some identifier. They both perform the same thing and are shorthand ways of referring to columns by name. You can use either all the above references when you read different people's spark code, I leave it to the reader for you to use whatever is most comfortable for you.

## Explicit Column References

If you need to refer to a specific DataFrame's column, you can use the `col` method on the specific DataFrame. This can be useful when you are performing a join and need to refer to a specific column in one DataFrame by its column name with another column in the joined DataFrame. We will see how to do joins chapter. As an added benefit, Spark does not need to resolve the column names itself (during the *analyzer* phase) because we did that for Spark.

```
df.col("count")
```

Now we mentioned that columns are expressions, so what is an expression? An expression is a set of transformations on one or more values in a DataFrame. Think of it like a function that takes as input one or more column names, resolves them and then potentially applies more expressions to produce a single value for each record in the dataset. Importantly, this “single value” can actually be a complex type like a Map type or Array type.

In the simplest case, an expression, created via the `expr` function, is a DataFrame column reference.

```
import org.apache.spark.sql.functions.{expr, col}
```

In this simple instance, `expr("someCol")` is equivalent to `col("someCol")`.

## Columns as Expressions

Columns provide a subset of expression functionality. If you use `expr` to wish to perform transformations on that column, you must perform a column reference. When using an expression, the `expr` function will parse transformations and column references from a string and subsequently be passed into further transformations. Let's look at some examples.

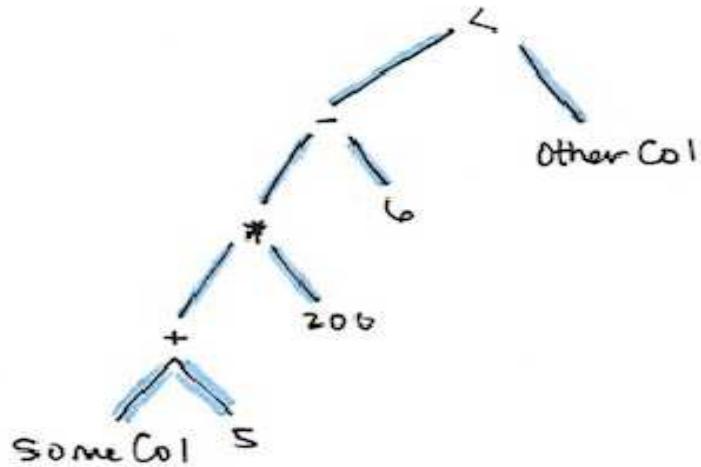
`expr("someCol - 5")` is the same transformation as performing `col("someCol") - 5` or even `expr("someCol") - 5`. That's because `expr` compiles these to a logical tree specifying the order of operations. This can be a bit confusing at first, but remember a couple of key points.

1. Columns are just expressions.
2. Columns and transformations of those column compile to the execution plan as parsed expressions.

Let's ground this with an example.

```
((col("someCol") + 5) * 200) - 6) < col("otherCol")
```

Figure 1 shows an illustration of that logical tree.



This might look familiar because it's a directed acyclic graph. It can be represented equivalently with the following code.

```
%scala  
  
import org.apache.spark.sql.functions.expr  
  
expr("(((someCol + 5) * 200) - 6) < otherCol")  
  
%python  
  
from pyspark.sql.functions import expr  
  
expr("(((someCol + 5) * 200) - 6) < otherCol")
```

This is an extremely important point to reinforce. Notice how the expression is actually valid SQL code as well, just like you might write in a SELECT statement? That's because this SQL expression and the DataFrame code compile to the same underlying logical tree prior to execution. This means you can write your expressions as DataFrame code as SQL expressions and get the exact same benefits. You likely saw this in the first chapters of the book and we covered this more extensively in the `myRow[U]` section.

```
    .lit(1).alias("One")  
).show(2)
```

In SQL, literals are just the specific value.

```
%sql  
  
SELECT  
    *,  
    1 as One  
FROM  
    dfTable  
LIMIT 2
```

This will come up when you might need to check if a date is greater than or equal to some value.

```
expr("some value")  
    .lit(1).as("something")  
).show(2)
```

```
%python
```

```
from pyspark.sql import expr, col, column  
  
df.select("select count(*) as NAME").show(2)  
.load("/mnt/defg/flight-data/json/2015-summary.json")  
df.createOrReplaceTempView("dfTable")
```

We can also create DataFrames on the fly by taking a set of rows and converting them to a DataFrame.

```
%scala  
  
import org.apache.spark.sql.Row  
import org.apache.spark.sql.types.{StructField, StructType, StringType, LongType}  
  
val myManualSchema = new StructType(Array(  
    new StructField("some", StringType, true),  
    new StructField("col", StringType, true),  
    new StructField("names", LongType, false) // just for fun  
)
```

11

```
val myRows = Seq(Row("Hello", null, 1L))
val myRDD = spark.sparkContext.parallelize(myRows)

val myDf = spark.createDataFrame(myRDD, myManualSchema)
myDf.show()
```



















