

```
In [1]: println("Hello World!")
```

```
Hello World!
```

```
In [2]: new java.io.File("/home/jovyan/work/data/shakespeare").list
```

```
Out[2]: Array(merrywivesofwindsor, twelfthnight, midsummersnightsdream, loveslabo  
urslost, asyoulikeit, comedyoferrors, muchadoaboutnothing, tamingoftheshr  
ew)
```

## Just Enough Scala for Spark

Dean Wampler, Ph.D. [@deanwampler](https://twitter.com/deanwampler) (<https://twitter.com/deanwampler>) ([email \(mailto:deanwampler@gmail.com\)](mailto:deanwampler@gmail.com))

Welcome. This notebook teaches you the core concepts of [Scala \(https://scala-lang.org\)](https://scala-lang.org) necessary to use [Apache Spark's \(https://spark.apache.org\)](https://spark.apache.org) Scala API effectively. Spark does a nice job exploiting the nicest features of Scala, while avoiding most of the more difficult and obscure features.

## Introduction: Why Scala?

Spark lets you use Scala, Java, Python, R, and SQL to do your work. Scala and Java appeal to *data engineers*, who do the heavy lifting of building resilient and scalable infrastructures for *Big Data*. Python, R, and SQL appeal to *data scientists*, who build models for analyzing data, including machine learning, as well as explore data interactively, where SQL is very convenient.

These aren't hard boundaries. Many people do both roles. Many data engineers like Python and may use SQL and R. Many data scientists have decided to use Scala with Spark.

Briefly, some of the advantages of using Scala include the following:

- **Performance:** Since Spark is written in Scala, you get the best performance and the most complete API coverage when you use Scala. It's true that with [DataFrames \(http://spark.apache.org/docs/latest/sql-programming-guide.html\)](http://spark.apache.org/docs/latest/sql-programming-guide.html), code written in all five languages performs about the same. If you need to use the [RDD \(http://spark.apache.org/docs/latest/programming-guide.html#resilient-distributed-datasets-rdds\)](http://spark.apache.org/docs/latest/programming-guide.html#resilient-distributed-datasets-rdds) API, then Scala provides the best performance, with Java a close second.
- **Debugging:** When runtime problems occur, understanding the exception stack frames and other debug information is easiest if you know Scala. Unfortunately, the "abstraction leaks" when problems occur.
- **Concise, Expressive Code:** Compared to Java, Scala code is much more concise and several features of Scala make your code even more concise. This elevates your productivity and makes it easier to imagine a design approach and then write it down without having to translate the idea to a less flexible API that reflects idiomatic language constraints. (You'll see this in action as we go.)
- **Type Safety:** Compared to Python and R, Scala code benefits from *static typing* with *type inference*. *Static typing* means that the Scala parser finds more errors in your expressions at

compile time, when they don't match expected types, rather than discovering the problem later at run time. However, *type inference* means you don't have to add a lot of explicit type information to your code. In most cases, Scala will infer the correct types for you.

## Why Not Scala?

Scala isn't perfect. There are two disadvantages compared to Python and R:

- **Libraries:** Python and R have a rich ecosystem of data analytics libraries. While the picture is improving for Scala, Python and R are still well ahead.
- **Advanced Language Features:** Mastering advanced language features gives you a lot of power to exploit, but if you don't understand those features, they can get in your way when you're just trying to get work done. Scala has some sophisticated constructs, especially in its *type system*. Fortunately, Spark mostly hides the advanced constructs.

## For More on Scala

I can only scratch the surface of Scala here. We'll "sketch" concepts without too much depth. You'll learn enough to make use of them, but eventually you'll want to deepen your understanding.

When you need more information, consider these resources:

- [Programming Scala, Second Edition](http://shop.oreilly.com/product/0636920033073.do) (<http://shop.oreilly.com/product/0636920033073.do>): My comprehensive introduction to Scala.
- [Scala Language Website](http://scala-lang.org/) (<http://scala-lang.org/>): Where to download Scala, find documentation (e.g., the [Scaladocs](http://www.scala-lang.org/api/current/#package) (<http://www.scala-lang.org/api/current/#package>): Scala library documentation, like [Javadocs](https://docs.oracle.com/javase/8/docs/api/) (<https://docs.oracle.com/javase/8/docs/api/>)), and other information.
- [Lightbend Scala Services](http://www.lightbend.com/services/) (<http://www.lightbend.com/services/>): training, consulting, and support for Scala.

## For More on Spark

- [Apache Spark website](https://spark.apache.org/) (<https://spark.apache.org/>): Spark downloads and documentation.
- [Databricks Spark Resources](https://sparkhub.databricks.com/resources/) (<https://sparkhub.databricks.com/resources/>): Databricks, the company founded by the creators of Sparks, offers free training materials. Consider the [Databricks Unified Data Analytics Platform](https://databricks.com/product/unified-data-analytics-platform) (<https://databricks.com/product/unified-data-analytics-platform>) for a full-featured, hosted data system including Spark.

For now, I recommend that you open the Scaladocs for Scala and for Spark's Scala API. Clicking these two links will open them in new browser tabs:

- Scaladocs for [Scala](http://www.scala-lang.org/api/current/#package) (<http://www.scala-lang.org/api/current/#package>).
- Scaladocs for [Spark](http://spark.apache.org/docs/latest/api/scala/index.html#package) (<http://spark.apache.org/docs/latest/api/scala/index.html#package>).

- Use the search bar in the upper-left-hand side to find a particular *type*. (For example, try "RDD" in the Spark Scaladocs.)
- To search for a particular *method*, click the character under the search box for the method name's first letter, then scroll to it.

## Prerequisites

I'll assume some prior programming experience in any language. Some familiarity with Java is assumed, but if you don't know Java, you should be able to search for explanations for anything unfamiliar.

This isn't an introduction to Spark itself. Some prior exposure to Spark is helpful, but I'll briefly explain most Spark concepts we'll encounter, too.

Throughout, you'll find links to more information on important topics.

## About Notebooks

You're using the [Jupyter](http://jupyter.org/) (<http://jupyter.org/>). [All Spark Notebook Docker image](https://hub.docker.com/r/jupyter/all-spark-notebook/) (<https://hub.docker.com/r/jupyter/all-spark-notebook/>). As described in the [GitHub README](https://github.com/deanwampler/JustEnoughScalaForSpark) (<https://github.com/deanwampler/JustEnoughScalaForSpark>) you import this notebook into Jupyter running as a Docker image.

Notebooks let you mix documentation, like this [Markdown](https://daringfireball.net/projects/markdown/) (<https://daringfireball.net/projects/markdown/>) "cell", with cells that contain code, graphs of results, etc. The metaphor is a physical notebook a scientist or student might use while working in a laboratory.

The menus and toolbar at the top provide options for evaluating a cell, adding and deleting cells, etc. You'll want to learn keyboard shortcuts if you use notebooks a lot.

### Tips:

1. Invoke the *Help > Keyboard Shortcuts* menu item, then capture the page as an image (it's a modal dialog, unfortunately). Learn a few shortcuts each day.
2. For now, just know that you can click into any cell to move the focus. When you're in a cell, `shift+enter` evaluates the cell (parses and renders the Markdown or runs the code), then moves to the next cell. Try it for a few cells. I'll wait...

Okay. It's particularly nice that you can edit a cell you've already evaluated and rerun it. This is great when you're experimenting with code.

## The Environment

Let's configure the environment to always show us the types of expressions.

```
In [3]: %showTypes on
```

Types will be printed.

When you start this notebook, the Jupyter Spark plugin creates a [SparkContext](http://spark.apache.org/docs/latest/programming-guide.html#initializing-spark) (<http://spark.apache.org/docs/latest/programming-guide.html#initializing-spark>) for you. This is the entry point of any Spark application (even when you use the newer `SparkSession`). It knows how to connect to your cluster (or run locally in the same JVM), how to configure properties, etc. It also runs a Web UI that lets you monitor your running jobs. The instance of `SparkContext` is called `sc`. The next cell simply confirms that it exists.

```
In [4]: sc
```

```
Out[4]: org.apache.spark.SparkContext = org.apache.spark.SparkContext@5fcfb180
```

Here are few useful bits of information:

```
In [5]: println("Spark version:      " + sc.version)
println("Spark master:      " + sc.master)
println("Running 'locally'?:" + sc.isLocal)
```

```
Spark version:      2.4.5
Spark master:      local[*]
Running 'locally'?: true
```

## Let's Load Some Data (and Start Learning Scala)

We're going to write real Spark programs and use them as vehicles for learning Scala and how to use it with Spark.

But first, we need to set up some text files we'll use, which contain some of the plays of Shakespeare. The next few cells define some helper methods (functions) to do this and then perform the steps. We'll start learning Scala concepts as we go.

### **Note:** "method" vs. "function"

Scala follows a common object-oriented convention where the term *method* is used for a function that's attached to a class or instance. Unlike Java, at least before Java 8, Scala also has *functions* that are not associated with a particular class or instance.

In our next code example, we'll define a few helper *methods* for printing information, but you won't see a class definition here. So, what class is associated with these methods? When you use Scala in a notebook, you're actually using the

Scala interpreter, which wraps any expressions and definitions we write into a hidden, generated class. The interpreter has to do this in order to generate valid JVM byte code.

Unfortunately, it can be a bit confusing when to use a method vs. a function, reflecting Scala's hybrid nature as an object-oriented and a functional language. Fortunately, in many cases, we can use methods and functions interchangeably, so we won't worry about the distinction too much from now on.

We're defining methods now. We'll see what a real *function* looks like soon.

Okay, here are two convenience methods for printing either an error message or a simple "information" message. We'll explain the syntax in a subsequent cell below.

```
In [6]: /*
  * "info" takes a single String argument, prints it on a line,
  * and returns it.
  */
def info(message: String): String = {
  println(message)

  // The last expression in the block, message, is the return value.
  // "return" keyword not required.
  // Do no additional formatting for the return string.
  message
}
```

info: (message: String)String

```
In [7]: /*
  * "error" takes a single String argument, prints a formatted error message
  * and returns the message.
  */
def error(message: String): String = {

  // Print the string passed to "println" and add a linefeed ("ln"):
  // See the next cell for an explanation of how the string is constructed
  val fullMessage = s"""
    | *****
    | ERROR: $message
    | *****
    | """.stripMargin
  println(fullMessage)

  fullMessage
}
```

error: (message: String)String

Let's try them.

```
In [8]: val infoString = info("All is well.")
```

```
All is well.
```

```
infoString: String = All is well.
```

```
Out[8]: infoString: String = All is well.
```

```
In [9]: val errorString = error("Uh oh...")
```

```
*****
```

```
ERROR: Uh oh...
```

```
*****
```

```
errorString: String =
```

```
"
```

```
*****
```

```
ERROR: Uh oh...
```

```
*****
```

```
"
```

```
Out[9]: errorString: String =
```

```
In [10]: errorString
```

```
"
```

```
*****
```

```
ERROR: Uh oh...
```

```
*****
```

```
"
```

```
Out[10]: String =
```

Method definitions have the following elements, in order:

- The `def` keyword.
- The method's name (`error` and `info` here).
- The argument list in parentheses. If there are no arguments, the empty parentheses can be omitted. This is common for `toString` and "getter"-like methods that simply return a field in an instance, etc.
- A colon followed by the type of the value returned by the method. This can often be inferred by Scala, so it's optional, but recommended for readability by users!
- An `=` (equals) sign that separates the method *signature* from the *body*.
- The body in braces `{ ... }`, although if the body consists of a single expression, the braces are optional.
- The last expression in the body is used as the return value. The `return` keyword is optional and rarely used.
- Semicolons (`;`) are inferred at the end of lines (in most cases) and rarely used.

Look at the argument list for `error`. It is `(message: String)`, where `message` is the argument name and its type is `String`. This convention for *type annotations*, `name: Type`, is also used for the return type, `error(...): String`. Type annotations are required by Scala for method arguments. They are optional in most cases for the return type. We'll see that Scala can infer the types of many expressions and variable declarations.

Scala uses the same comment conventions as Java, `// ...` for a single line, and `/* ... */` for a comment block.

#### Note: Expression vs. Statement

An *expression* has a value, while a *statement* does not. Hence, when we assign an expression to a variable, the value the expression returns is assigned to the variable.

Inside `error`, we used a combination *interpolated* and *triple-quoted* string with the syntax `s"""..."""`:

- **Triple-quoted string:** `"""..."""`. Useful for embedding newlines, like we did inside `error`. (We'll see another benefit later.)
- **String interpolation:** Invoked by putting `s` before the string, e.g., `s"..."` or `s"""..."""`. Lets us embed variable references and expressions, where the string conversion will be inserted automatically. For example:

```
In [11]: s"""Use braces for expressions: ${sc.version}.
You can omit the braces when just using a variable: $sc
However, watch for ambiguities like ${sc}andextrastuff"""
```

```
Use braces for expressions: 2.4.5.
You can omit the braces when just using a variable: org.apache.spark.SparkContext@5fcfb180
However, watch for ambiguities like org.apache.spark.SparkContext@5fcfb180andextrastuff
```

```
Out[11]: String =
```

Another feature we're using for triple-quoted strings is the ability to strip the leading whitespace off each line. The `stripMargin` method removes all whitespace before and including the `|`. This lets you indent those lines for proper code formatting, but not have that whitespace remain in the string. In the following example, the resulting string has blank lines at the beginning and end. Note what happens with whitespace before `line2` and `line3` when the full string is printed:

```
In [12]: s"""
          |line 1
          |  line 2
          ||  line 3
          |""".stripMargin
```

```
"
line 1
  line 2
    line 3
"
```

```
Out[12]: String =
```

Character "literals" are specified single quotes, `'/'`, while strings use double quotes, `"/"`.

```
In [13]: '/'
```

```
Out[13]: Char = /
```

```
In [14]: "/"
```

```
Out[14]: String = /
```

## Mutable Variables vs. Immutable Values

See how how to declare an immutable value before with `val`. Let's explore this a bit more:

- `val immutableValue = ...`: Once initialized, we can't assign a *different* value to `immutableValue`.
- `var mutableVariable = ...`: We can assign new values to `mutableVariable` as often as we want.

It's *highly recommended* that you only use `vals` unless you have a good reason for needing mutability, which is a very common source of bugs!!

A `val immutableValue` could point to an instance that itself *is* mutable, e.g., an [Array](http://www.scala-lang.org/api/current/#scala.Array) (<http://www.scala-lang.org/api/current/#scala.Array>) (Scala uses Java arrays, which are mutable). In this case, while we can't assign a new array to `immutableValue`, we can change elements within the array! Put another way, immutability isn't *transitive*.

## Setup the Files

The notebook already has the data files we need, several of Shakespeare's plays. They are in the `/home/jovyan/work/data/shakespeare` subdirectory in the container ( `data/shakespeare` in the git project). There is one file for each play.

We'll write some Scala code to verify they are there, primarily so we can learn some more Scala.



Many of the types used in Scala code are from Java's library (JDK). Because Scala compiles to JVM byte code, you can use any Java library you want from Scala. We've been using [java.lang.String](https://docs.oracle.com/javase/8/docs/api/java/lang/String.html) (<https://docs.oracle.com/javase/8/docs/api/java/lang/String.html>). Now we'll use [java.io.File](https://docs.oracle.com/javase/8/docs/api/java/io/File.html) (<https://docs.oracle.com/javase/8/docs/api/java/io/File.html>) to work with files and directories. As before, we'll use comments to explain a few other new Scala constructs.

```
In [15]: // Import File. Unlike Java, the semicolon ';' is not required.  
import java.io.File
```

Here the the directory where the files should be located.

```
In [16]: val shakespeare = new File("/home/jovyan/work/data/shakespeare")
```

```
shakespeare: java.io.File = /home/jovyan/work/data/shakespeare
```

```
Out[16]: shakespeare: java.io.File = /home/jovyan/work/data/shakespeare
```

Scala's `if` construct is actually an expression (in Java they are `_statements_`). The `if` expression will return `true` or `false` and assign it to `success`, which we'll use in a moment.

```
In [17]: val success = if (shakespeare.exists == false) { // doesn't exist already  
    error(s"Data directory path doesn't exist! $shakespeare") // ignore re  
    false  
} else {  
    info(s"$shakespeare exists")  
    true  
}  
println("success = " + success)
```

```
/home/jovyan/work/data/shakespeare exists
```

```
success = true
```

```
success: Boolean = true
```

```
Out[17]: success: Boolean = true
```

Now lets verify the files we expect are there, again to learn some more Scala.

```
In [18]: val pathSeparator = File.separator
val targetDirName = shakespeare.toString
val plays = Seq(
  "tamingoftheshrew", "comedyoferrors", "loveslabourslost", "midsummersni
  "merrywivesofwindsor", "muchadoaboutnothing", "asyoulikeit", "twelfthni

if (success) {
  println(s"Checking that the plays are in $shakespeare:")
  val failures = for {
    play <- plays
    playFileName = targetDirName + pathSeparator + play
    playFile = new File(playFileName)
    if (playFile.exists == false)
  } yield {
    s"$playFileName:\tNOT FOUND!"
  }

  println("Finished!")
  if (failures.size == 0) {
    info("All plays found!")
  } else {
    println("The following expected plays were not found:")
    failures.foreach(play => error(play))
  }
}
```

Checking that the plays are in /home/jovyan/work/data/shakespeare:

Finished!

All plays found!

```
pathSeparator: String = /
targetDirName: String = /home/jovyan/work/data/shakespeare
plays: Seq[String] = List(tamingoftheshrew, comedyoferrors, loveslaboursl
ost, midsummersnightsdream, merrywivesofwindsor, muchadoaboutnothing, asy
oulikeit, twelfthnight)
```

Out[18]: Any = All plays found!

I'm using a so-called *for comprehension*. They are *expressions*, not *statements* like Java's *for* loops. They have the form:

```
for {
  play <- plays
  ...
} yield { block_of_final_expressions }
```

We iterate through a collection, `plays`, and assign each one to the `play` variable (actually an immutable value for each pass through the loop).

After assigning to `play`, subsequent steps in the *for* comprehension use it. First, a [java.io.File](https://docs.oracle.com/javase/8/docs/api/java/io/File.html) (<https://docs.oracle.com/javase/8/docs/api/java/io/File.html>) instance, `playFile`, is created. Then, `playFile` is used to evaluate a conditional - does the file already exist? (It should!)

If the file already exists, the conditional returns `false`, which short-circuits the loop and goes to the next `play` in the list. If the file doesn't exist, the `yield` keyword tells Scala that I want to use the expression that follows to construct a new element, an *interpolated* string, for the missing play. From those returned elements, zero or more, a new collection is constructed. The final `if` block determines if the new collection has zero elements (expected), then prints an `info` message. If there were missing files, an `error` message is printed for each one of them.

## Passing Functions as Arguments

Note how we printed the returned `successes` collection of strings. The idiom `collection.foreach(println)` is handy for looping over the elements and printing them, one per line. But how exactly does this work? (We'll use `plays` instead of `failures`, because the latter should be empty!)

```
In [19]: println("Pass println as the function to use for each element:")
plays.foreach(println)

println("\nUsing an anonymous function that calls println: `str => println(str)`")
println("(Note that the type of the argument `str` is inferred to be String)")
plays.foreach(str => println(str))

println("\nAdding the argument type explicitly. Note that the parentheses are required.")
plays.foreach((str: String) => println(str))

println("\nWhy do we need to name this argument? Scala lets us use `_` as a placeholder.")
plays.foreach(println(_))

println("\nFor longer functions, you can use `{...}` instead of `(...)`.")
println("Why? Because it gives you the familiar multiline block syntax with curly braces.")
plays.foreach {
  (str: String) => println(str)
}

println("\nThe `_` placeholder can be used *once* for each argument in the lambda.")
println("As an example, use `reduceLeft` to sum some integers.")
val integers = 0 to 10 // Return a "range" from 0 to 10, inclusive
integers.reduceLeft((i,j) => i+j)
integers.reduceLeft(_+_)
```

Pass println as the function to use for each element:

```
tamingoftheshrew
comedyoferrors
loveslabourslost
midsummersnightsdream
merrywivesofwindsor
muchadoaboutnothing
asyoulikeit
twelfthnight
```

Using an anonymous function that calls println: `str => println(str)`  
(Note that the type of the argument `str` is inferred to be String.)

```
tamingoftheshrew
comedyoferrors
loveslabourslost
midsummersnightsdream
merrywivesofwindsor
muchadoaboutnothing
asyoulikeit
twelfthnight
```

Adding the argument type explicitly. Note that the parentheses are required.

```
tamingoftheshrew
comedyoferrors
loveslabourslost
midsummersnightsdream
merrywivesofwindsor
muchadoaboutnothing
asyoulikeit
twelfthnight
```

Why do we need to name this argument? Scala lets us use `_` as a placeholder.

```
tamingoftheshrew
comedyoferrors
loveslabourslost
midsummersnightsdream
merrywivesofwindsor
muchadoaboutnothing
asyoulikeit
twelfthnight
```

For longer functions, you can use `{...}` instead of `(...)`.

Why? Because it gives you the familiar multiline block syntax with `{...}`

```
tamingoftheshrew
comedyoferrors
loveslabourslost
midsummersnightsdream
merrywivesofwindsor
muchadoaboutnothing
asyoulikeit
twelfthnight
```

The `_` placeholder can be used *once* for each argument in the list.

As an assume, use ``reduceLeft`` to sum some integers.

```
integers: scala.collection.immutable.Range.Inclusive = Range(0, 1, 2, 3,
4, 5, 6, 7, 8, 9, 10)
```

Out[19]: Int = 55

## Our First Spark Program

Whew! We've learned a lot of Scala already while doing typical data science chores (i.e., fetching data). Now let's implement a real algorithm using Spark, *Inverted Index*.

## Inverted Index - When You're Tired of Counting Words...

You'll want use *Inverted Index* when you create your next "Google killer". It takes in a corpus of documents (e.g., web pages), tokenizes the words, and outputs for each word a list of the documents that contain it, along with the corresponding counts.

This is a slightly more interesting algorithm than `_Word Count_`, the classic "hello world" program everyone implements when they learn Spark.

The term *inverted* here means we start with the words as part of the input *values*, while the *keys* are the document identifiers, and we'll switch ("invert") to using the words as keys and the document identifiers as values.

Here's our first version, all at once. This is *one, long expression*. Note the periods `.` at the end of the subexpressions.

```
In [20]: val iiFirstPass1 = sc.wholeTextFiles(shakespeare.toString).
    flatMap { location_contents_tuple2 =>
      val words = location_contents_tuple2._2.split(" "\W+"")
      val fileName = location_contents_tuple2._1.split(pathSeparator).last
      words.map(word => ((word, fileName), 1))
    }.
    reduceByKey((count1, count2) => count1 + count2).
    map { word_file_count_tup3 =>
      (word_file_count_tup3._1._1, (word_file_count_tup3._1._2, word_file
    }.
    groupByKey.
    sortByKey(ascending = true).
    mapValues { iterable =>
      val vect = iterable.toVector.sortBy { file_count_tup2 =>
        (-file_count_tup2._2, file_count_tup2._1)
      }
      vect.mkString(",")
    }
  }
```

```
iiFirstPass1: org.apache.spark.rdd.RDD[(String, String)] = MapPartitionsRDD[9] at mapValues at <console>:44
```

```
Out[20]: iiFirstPass1: org.apache.spark.rdd.RDD[(String, String)] = MapPartitionsRDD[9] at mapValues at <console>:44
```

Now let's break it down into steps, assigning each step to a variable. This extra verbosity let's us see what Scala infers for the type returned by each expression, helping us learn.

This is one of the nice features of Scala. We don't have to put in the type information ourselves, most of the time, like we would have to do for Java code. Instead, we let the compiler give us feedback about what we just created. This is especially useful when you're learning a new API, like Spark's.

```
In [21]: val fileContents = sc.wholeTextFiles(shakespeare.toString)
    fileContents // force the notebook to print the type.
```

```
fileContents: org.apache.spark.rdd.RDD[(String, String)] = /home/jovyan/work/data/shakespeare MapPartitionsRDD[11] at wholeTextFiles at <console>:30
```

```
Out[21]: org.apache.spark.rdd.RDD[(String, String)] = /home/jovyan/work/data/shakespeare MapPartitionsRDD[11] at wholeTextFiles at <console>:30
```

The second line, with `fileContents` by itself, is there so the notebook will show us its type information. (Try to remove it and re-evaluate the cell. Nothing is printed.).

The output is telling us that `fileContents` has the type `RDD[(String,String)]` (<http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.rdd.RDD>), but `RDD` is a base class and the actual instance is a `MapPartitionsRDD`, which is a "private" implementation subclass of `RDD`.

A name followed by square brackets, `[...]`, means that `RDD[...]` requires one or more type parameters in the brackets. In this case, a single type parameter, which represents the type of the records held by the `RDD`.

The single type parameter is given by `(String,String)`, which is a convenient shorthand for `Tuple2[String,String]` (<http://www.scala-lang.org/api/current/index.html#scala.Tuple2>). That is, we have two-element *tuples* as records, where the first element is a `String` for a file's fully-qualified path and the second element is a `String` for the contents of that file. This is what `SparkContext.wholeTextFiles` returns for us. We'll use the path to remember where we found words, while the contents contains the words themselves (of course).

To recap, the following two types are equivalent:

- `RDD[(String,String)]` - Note parentheses nested in brackets, `[ (... ) ]`.
- `RDD[Tuple2[String,String]]` - Note nested brackets `[... [...]]`, not `[ (... ) ]`.

We'll see shortly that you can also write *instances* of `Tuple2[T1,T2]` (<http://www.scala-lang.org/api/current/index.html#scala.Tuple2>) with the same syntax, e.g., `("foo", 101)`, for a `(String,Int)` tuple, and similarly for *higher-arity* tuples (up to 22 elements...), e.g., `("foo", 101, 3.14159, ("bar", 202L))`. Run the next cell to see the type signature for this last tuple.

```
In [22]: ("foo", 101, 3.14159, ("bar", 202L))
```

```
Out[22]: (String, Int, Double, (String, Long)) = (foo,101,3.14159,(bar,202))
```

Do you understand it? Do you see that it's a four-element tuple and not a five-element tuple? This is because the `("bar", 202L)` is a nested tuple. It's the fourth element of the outer tuple.

**Exercise:** Try creating some more tuples with elements of different types. Use the next cell.

```
In [23]: (1,2)
```

```
Out[23]: (Int, Int) = (1,2)
```

How many `fileContents` records do we have? Not many. It should be the same number as the number of files we downloaded above.

```
In [24]: fileContents.count
```

```
Out[24]: Long = 8
```

**NOTE:** We called the `RDD.count` method, whereas most Scala collections have a `size` method.

Now for our next step in the calculation. First, "tokenize" the contents into words by splitting on

non-alphanumeric characters, meaning all runs of whitespace (including the newlines), punctuation, etc.

Next, the fully-qualified path is verbose and the same prefix is repeated for all the files, so let's extract just the last element of it, the unique file name.

Then form new tuples with the words and file names.

**Note:** This "tokenization" approach is very crude. It improperly handles contractions, like *it's* and hyphenated words like *world-changing*. When you kill Google, be sure to use a real *natural language processing* (NLP) tokenization technique.

```
In [25]: val wordFileNameOnes = fileContents.flatMap { location_contents_tuple2 =>
// example input record: (file_path, "all the words in the file")
// mytuple._2 => give me the 2nd element
val words = location_contents_tuple2._2.split(" "\\W+" ")
// mytuple._1 => give me the 1st element
val fileName = location_contents_tuple2._1.split(pathSeparator).last
// create a new tuple to return. Note how we structured it!
words.map(word => ((word, fileName), 1))
}
wordFileNameOnes
```

```
wordFileNameOnes: org.apache.spark.rdd.RDD[((String, String), Int)] = MapPartitionsRDD[12] at flatMap at <console>:30
```

```
Out[25]: org.apache.spark.rdd.RDD[((String, String), Int)] = MapPartitionsRDD[12] at flatMap at <console>:30
```

I find this hard to read and shortly I'll show you a much more elegant, alternative syntax.

Let's understand the difference between `map` and `flatMap`. If I called `fileContents.map`, it would return exactly *one* new record for each record in `fileContents`. What I actually want instead are new records for each word-fileName combination, a significantly larger number (but the data in each record will be much smaller).

Using `fileContents.flatMap` gives me what I want. Instead of returning one output record for each input record, a `flatMap` returns a *collection* of new records, zero or more, for *each* input record. These collections are then *flattened* into one big collection, another `RDD` in this case.

What should `flatMap` actually do with each record? I pass a *function* to define what to do. I'm using an unnamed or *anonymous* function. The syntax is `argument_list => body`:

```
location_contents_tuple2 =>
  val words = ...
  ...
}
```



I have a single argument, the record, which I named `location_contents_tuple2`, a verbose way to say that it's a two-element tuple with an input file's location and contents. I don't require a type parameter after `location_contents_tuple2`, because it's inferred by Scala. The `=>` "arrow" separates the argument list from the body, which appears on the next few lines.

When a function takes more than one argument or you add explicit type *annotations* (e.g., `: (String, Int, Double)`), then you need parentheses. Here are three examples:

```
(some_tuple3: (String, Int, Double)) => ...
(arg1, arg2, arg3) => ...
(arg1: String, arg2: Int, arg3: Double) => ...
```

We're letting Scala infer the argument type in our case, `(String, String)`.

Wait, I said we're passing a function as an argument to `flatMap`. If so, why am I using braces `{...}` around this function argument instead of parentheses `(...)` like you would normally expect when passing arguments to a method like `flatMap`?

It's because Scala lets us substitute braces instead of parentheses so we have the familiar block-like syntax `{...}` we know and love for `if` and `for` expressions. I could use either braces or parentheses here. The convention in the Scala community is to use braces for a multi-line anonymous function and to use parentheses for a single expression when it fits on the same line.

Now, for each `location_contents_tuple2`, I access the *first* element using the `_1` method and the *second* element using `_2`.

The file `contents` is in the second element. I split it by calling Java's `String.split` method, which takes a *regular expression* string. Here I specify a regular expression for one or more, non-alphanumeric characters. `String.split` returns an `Array[String]` of the words.

```
val words = location_contents_tuple2._2.split(" "\\W+")
```

For the first tuple element, I extract the file name at the end of the location path. This isn't necessary, but it makes the output more readable if I remove the long, common prefix from the path.

```
val fileName = location_contents_tuple2._1.split(pathSeparator).last
```

Finally, still inside the anonymous function passed to `flatMap`, I use Scala's `Array.map` (*not* `RDD.map`) to transform each `word` into a tuple of the form `((word, fileName), 1)`.

```
words.map(word => ((word, fileName), 1))
```

Why did I embed a tuple of `(word, fileName)` inside the "outer" tuple with a `1` as the second element? Why not just write a three-element tuple, `(word, fileName, 1)`? It's because I'll use the `(word, fileName)` as a *key* in the next step, where I'll find all unique word-fileName combinations (using the equivalent of a `group by` statement). So, using the nested `(word, fileName)` as my *key* is most convenient. The `1` *value* is a "seed" count, which I'll use to count the occurrences of the unique `(word, fileName)` pairs.

**Notes:**

- For historical reasons, tuple indices start at 1, not 0. Arrays and other Scala collections index from 0.
- I said previously that *method* arguments have to be declared with types. That's usually *not* required for *function* arguments, as here.
- Another benefit of triple-quoted strings that makes them nice for regular expressions is that you don't have to escape regular expression metacharacters, like `\w`. If I used a single-quoted string, I would have to write it as `"\\w+"`. Your choice...

Let's count the number of records we have and look at a few of the lines. We'll use the `RDD.take` method to grab the first 10 lines, then loop over them and print them.

```
In [26]: wordFileNameOnes.count
```

```
Out[26]: Long = 173336
```

```
In [27]: wordFileNameOnes.take(10).foreach(println)
```

```
((,merrywivesofwindsor),1)
((THE,merrywivesofwindsor),1)
((MERRY,merrywivesofwindsor),1)
((WIVES,merrywivesofwindsor),1)
((OF,merrywivesofwindsor),1)
((WINDSOR,merrywivesofwindsor),1)
((DRAMATIS,merrywivesofwindsor),1)
((PERSONAE,merrywivesofwindsor),1)
((SIR,merrywivesofwindsor),1)
((JOHN,merrywivesofwindsor),1)
```

We asked for results, so we forced Spark to run a job to compute results. Spark pipelines, like `iiFirstPass1` are `_lazy_`; nothing is computed until we ask for results.

When you're learning, it's useful to print some data to better understand what's happening. Just be aware of the extra overhead of running lots of Spark jobs.

The first record shown has "" (blank) as the word:

```
((,asyoulikeit),1)
```

Also, some words have all capital letters:

```
((DRAMATIS,asyoulikeit),1)
```

(You can see where these capitalized words occur if you look in the original files.) Later on, We'll filter out the blank-word records and use lower case for all words.

Now, let's join all the unique `(word,fileName)` pairs together.

```
In [28]: val uniques = wordFileNameOnes.reduceByKey((count1, count2) => count1 + count2)
uniques
```

```
uniques: org.apache.spark.rdd.RDD[(String, String), Int] = ShuffledRDD
[13] at reduceByKey at <console>:28
```

```
Out[28]: org.apache.spark.rdd.RDD[(String, String), Int] = ShuffledRDD[13] at re
duceByKey at <console>:28
```

In SQL you would use `GROUP BY` for this (including SQL queries you might write with Spark's [DataFrame](http://spark.apache.org/docs/latest/sql-programming-guide.html) (<http://spark.apache.org/docs/latest/sql-programming-guide.html>) API). However, in the `RDD` API, this is too expensive for our needs, because we don't care about the groups themselves, the long list of repeated `(word, fileName)` pairs. We only care about how many elements are in each group, that is their *size*. That's the purpose of the `1` in the tuples and the use of `RDD.reduceByKey`. It brings together all records with the same key, the unique `(word, fileName)` pairs, and then applies the anonymous function to "reduce" the values, the `1`s. I simply sum them up to compute the group counts.

Note that the anonymous function `reduceByKey` expects must take two arguments, so I need parentheses around the argument list. Since this function fits on the same line, I used parentheses for `reduceByKey`, instead of braces.

**Note:** All the `*ByKey` methods operate on two-element tuples and treat the first element as the key, by default.

How many are there? Let's see a few:

```
In [29]: uniques.count
```

```
Out[29]: Long = 27276
```

As you would expect from a `GROUP BY`-like statement, the number of records is smaller than before. There are about 1/6 as many records now, meaning that on average, each `(word, fileName)` combination appears 6 times.

```
In [30]: uniques.take(30).foreach(println)

((dexterity,merrywivesofwindsor),1)
((force,muchadoaboutnothing),2)
((whole,comedyoferrors),2)
((lamb,muchadoaboutnothing),2)
((blunt,tamingoftheshrew),3)
((letter,merrywivesofwindsor),19)
((crest,asyoulikeit),1)
((bestow,asyoulikeit),1)
((rear,midsummersnightsdream),1)
((crossing,tamingoftheshrew),1)
((wronged,merrywivesofwindsor),4)
((S,tamingoftheshrew),10)
((HIPPOLYTA,midsummersnightsdream),19)
((revolve,twelfthnight),1)
((er,merrywivesofwindsor),11)
((renown,asyoulikeit),1)
((cubiculo,twelfthnight),1)
((All,twelfthnight),3)
((power,loveslabourslost),8)
((Albeit,asyoulikeit),1)
((lips,tamingoftheshrew),3)
((upshot,twelfthnight),1)
((approach,midsummersnightsdream),4)
((mean,muchadoaboutnothing),5)
((embossed,asyoulikeit),1)
((varnish,loveslabourslost),2)
((Apollo,midsummersnightsdream),1)
((spangled,midsummersnightsdream),1)
((gentlemen,comedyoferrors),1)
((Rebuke,loveslabourslost),1)
```

For `_inverted index_`, we want our final keys to be the words themselves, so let's restructure the tuples from `((word,fileName),count)` to `(word,(fileName,count))`. Now, I'll still output two-element, key-value tuples, but the `word` will be the key and the `(fileName,count)` tuple will be the value.

```
In [31]: val words = uniques.map { word_file_count_tup3 =>
  (word_file_count_tup3._1._1, (word_file_count_tup3._1._2, word_file_cou
  })
```

```
words: org.apache.spark.rdd.RDD[(String, (String, Int))] = MapPartitionsRDD[14] at map at <console>:28
```

```
Out[31]: words: org.apache.spark.rdd.RDD[(String, (String, Int))] = MapPartitionsRDD[14] at map at <console>:28
```

The nested tuple methods, e.g., `_1._2`, are hard to read, making the logic somewhat obscure. We'll see a beautiful and elegant alternative shortly.

Now I'll use an actual `group by` operation, because I now need to retain the groups. Calling `RDD.groupByKey` uses the first tuple element, now just the `words`, to bring together all

occurrences of the unique words. Next, I'll sort the result by word, ascending alphabetically.

```
In [32]: val wordGroups = words.groupByKey.sortByKey(ascending = true)
wordGroups
```

```
wordGroups: org.apache.spark.rdd.RDD[(String, Iterable[(String, Int)])] =
ShuffledRDD[18] at sortByKey at <console>:28
```

```
Out[32]: org.apache.spark.rdd.RDD[(String, Iterable[(String, Int)])] = ShuffledRDD
[18] at sortByKey at <console>:28
```

Note that each group is actually a Scala [Iterable](http://www.scala-lang.org/api/current/index.html#scala.collection.Iterable) (<http://www.scala-lang.org/api/current/index.html#scala.collection.Iterable>), i.e., an abstraction for some sort of collection. (It's actually a Spark-defined, private collection type called a `CompactBuffer`.)

```
In [33]: wordGroups.count
```

```
Out[33]: Long = 11951
```

```
In [34]: wordGroups.take(30).foreach(println)
```

```
(,CompactBuffer((tamingoftheshrew,1), (asyoulikeit,1), (merrywivesofwindsor,1), (comedyoferrors,1), (midsummersnightsdream,1), (twelfthnight,1), (loveslabourslost,1), (muchadoaboutnothing,1)))
(A,CompactBuffer((loveslabourslost,78), (midsummersnightsdream,39), (muchadoaboutnothing,31), (merrywivesofwindsor,38), (comedyoferrors,42), (asyoulikeit,34), (twelfthnight,47), (tamingoftheshrew,59)))
(ABOUT,CompactBuffer((muchadoaboutnothing,18)))
(ACT,CompactBuffer((asyoulikeit,22), (comedyoferrors,11), (tamingoftheshrew,12), (loveslabourslost,9), (muchadoaboutnothing,17), (twelfthnight,18), (merrywivesofwindsor,23), (midsummersnightsdream,9)))
(ADAM,CompactBuffer((asyoulikeit,16)))
(ADO,CompactBuffer((muchadoaboutnothing,18)))
(ADRIANA,CompactBuffer((comedyoferrors,85)))
(ADRIANO,CompactBuffer((loveslabourslost,111)))
(AEGEON,CompactBuffer((comedyoferrors,20)))
(AEMELIA,CompactBuffer((comedyoferrors,16)))
(AEMILIA,CompactBuffer((comedyoferrors,3)))
(AEacides,CompactBuffer((tamingoftheshrew,1)))
(AEgeon,CompactBuffer((comedyoferrors,7)))
(AEgle,CompactBuffer((midsummersnightsdream,1)))
(AEmilia,CompactBuffer((comedyoferrors,4)))
(AEsculapius,CompactBuffer((merrywivesofwindsor,1)))
(AGUECHEEK,CompactBuffer((twelfthnight,2)))
(ALL,CompactBuffer((midsummersnightsdream,2), (tamingoftheshrew,2)))
(AMIENS,CompactBuffer((asyoulikeit,16)))
(ANDREW,CompactBuffer((twelfthnight,104)))
(ANGELO,CompactBuffer((comedyoferrors,36)))
(ANN,CompactBuffer((merrywivesofwindsor,1)))
(ANNE,CompactBuffer((merrywivesofwindsor,27)))
(ANTIPHOLUS,CompactBuffer((comedyoferrors,195)))
(ANTONIO,CompactBuffer((muchadoaboutnothing,32), (twelfthnight,32)))
(ARMADO,CompactBuffer((loveslabourslost,111)))
(AS,CompactBuffer((asyoulikeit,24)))
(AUDREY,CompactBuffer((asyoulikeit,18)))
(Abate,CompactBuffer((midsummersnightsdream,1), (loveslabourslost,1)))
(Abbess,CompactBuffer((comedyoferrors,2)))
```

Finally, let's clean up these `CompactBuffers`. Let's convert each to a Scala [Vector](http://www.scala-lang.org/api/current/index.html#scala.collection.immutable.Vector) (<http://www.scala-lang.org/api/current/index.html#scala.collection.immutable.Vector>) (a collection with  $O(1)$  performance for most operations), then sort it *descending* by count, so the locations that mention the corresponding word the *most* appear *first* in the list. (Think about how you would want a search tool to work...)

Note we're using `Vector.sortBy`, not an RDD sorting method. It takes a function that accepts each collection element and returns something used to sort the collection. By returning `(-fileNameCountTuple2._2, fileNameCountTuple2)`, I effectively say, "sort by the counts *descending* first, then sort by the file names." Why does `-fileNameCountTuple2._2` cause counts to be sorted descending, because I'm returning the negative of the value, so larger counts will be less than smaller counts, e.g., `-3 < -2`.

Finally, I take the resulting `Vector` and make a comma-separated string with the elements, using the helper method `mkString`.

What's `RDD.mapValues` ? I could use `RDD.map` , but I'm not changing the keys (the words), so rather than have to deal with the tuple with both elements, `mapValues` just passes in the value part of the tuple and reconstructs new `(key,value)` tuples with the new value that my function returns. So, `mapValues` is more convenient to use than `map` when I have two-element tuples and I'm not modifying the keys.

```
In [35]: val iiFirstPass2 = wordGroups.mapValues { iterable =>
          val vect = iterable.toVector.sortBy { file_count_tup2 =>
            (-file_count_tup2._2, file_count_tup2._1)
          }
          vect.mkString(",")
        }
```

```
iiFirstPass2: org.apache.spark.rdd.RDD[(String, String)] = MapPartitionsRDD[19] at mapValues at <console>:28
```

```
Out[35]: iiFirstPass2: org.apache.spark.rdd.RDD[(String, String)] = MapPartitionsRDD[19] at mapValues at <console>:28
```

We're done! The number of records is the same as for `wordGroups` (do you understand why?), so let's just see some of the records.

```
In [36]: iiFirstPass2.take(30).foreach(println)
```

```
(, (asyoulikeit,1), (comedyoferrors,1), (loveslabourslost,1), (merrywivesofwindsor,1), (midsummersnightsdream,1), (muchadoaboutnothing,1), (tamingoftheshrew,1), (twelfthnight,1))
(A, (loveslabourslost,78), (tamingoftheshrew,59), (twelfthnight,47), (comedyoferrors,42), (midsummersnightsdream,39), (merrywivesofwindsor,38), (asyoulikeit,34), (muchadoaboutnothing,31))
(ABOUT, (muchadoaboutnothing,18))
(ACT, (merrywivesofwindsor,23), (asyoulikeit,22), (twelfthnight,18), (muchadoaboutnothing,17), (tamingoftheshrew,12), (comedyoferrors,11), (loveslabourslost,9), (midsummersnightsdream,9))
(ADAM, (asyoulikeit,16))
(ADO, (muchadoaboutnothing,18))
(ADRIANA, (comedyoferrors,85))
(ADRIANO, (loveslabourslost,111))
(AEGEON, (comedyoferrors,20))
(AEMELIA, (comedyoferrors,16))
(AEMILIA, (comedyoferrors,3))
(AEAcides, (tamingoftheshrew,1))
(AEgeon, (comedyoferrors,7))
(AEgle, (midsummersnightsdream,1))
(AEmilia, (comedyoferrors,4))
(AEsculapius, (merrywivesofwindsor,1))
(AGUECHEEK, (twelfthnight,2))
(ALL, (midsummersnightsdream,2), (tamingoftheshrew,2))
(AMIENS, (asyoulikeit,16))
(ANDREW, (twelfthnight,104))
(ANGELO, (comedyoferrors,36))
(ANN, (merrywivesofwindsor,1))
(ANNE, (merrywivesofwindsor,27))
(ANTIPHOLUS, (comedyoferrors,195))
(ANTONIO, (muchadoaboutnothing,32), (twelfthnight,32))
(ARMADO, (loveslabourslost,111))
(AS, (asyoulikeit,24))
(AUDREY, (asyoulikeit,18))
(Abate, (loveslabourslost,1), (midsummersnightsdream,1))
(Abbess, (comedyoferrors,2))
```

Okay. Looks reasonable.

Next, I'll refine the code using a very powerful feature, `_pattern matching_`, which both makes the code more concise and easier to understand. It's my *favorite* feature of Scala.

Before I do that, try a few refinements on your own.

### Exercises:

- Add a filter statement to remove the first entry for the blank word `""`. You could do this one of two ways, using another "step" with [RDD.filter](http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.rdd.RDD) (<http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.rdd.RDD>) (search the [Scaladoc page] (<http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.rdd.RDD>) (<http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.rdd.RDD>)) for the `filter` method), or using the similar Scala collections method, [scala.collection.Seq.filter](http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.rdd.RDD)



(<http://www.scala-lang.org/api/current/index.html#scala.collection.Seq>). Both versions take a *predicate* function, one that returns `true` if the record should be *retained* and `false` otherwise. Do you think one choice is better than the other? Why? Or, are they basically the same? Reasons might include code comprehension and performance of one over the other.

- Convert all words to lower case. Calling `toLowerCase` on a string is all you need. Where's a good place to insert this logic?

I'll implement both changes in subsequent refinements below.

**NOTE:** If you would prefer to make a copy of the code in a new cell, use the *Insert* menu above to add cells. Or, learn another keyboard shortcut; `ESC` (escape key), followed by `A` for insert before or `B` for insert after. Then hit return to edit. Note the toolbar pop-down for setting the format of the cell. This cell you're reading is *Markdown*. Make sure to use *Code* for your source code cells.

## Pattern Matching

We've studied a real program and we've learned quite a bit of Scala. Let's improve it with my favorite Scala feature, *pattern matching*.

Here's the "first pass" version again for easy reference.

```
In [37]: val iiFirstPass1b = sc.wholeTextFiles(shakespeare.toString).
  flatMap { location_contents_tuple2 =>
    val words = location_contents_tuple2._2.split(" "\W+"")
    val fileName = location_contents_tuple2._1.split(pathSeparator).last
    words.map(word => ((word, fileName), 1))
  }.
  reduceByKey((count1, count2) => count1 + count2).
  map { word_file_count_tup3 =>
    (word_file_count_tup3._1._1, (word_file_count_tup3._1._2, word_file
  }.
  groupByKey.
  sortByKey(ascending = true).
  mapValues { iterable =>
    val vect = iterable.toVector.sortBy { file_count_tup2 =>
      (-file_count_tup2._2, file_count_tup2._1)
    }
    vect.mkString(", ")
  }
}
```

```
iiFirstPass1b: org.apache.spark.rdd.RDD[(String, String)] = MapPartitions
RDD[29] at mapValues at <console>:46
```

```
Out[37]: iiFirstPass1b: org.apache.spark.rdd.RDD[(String, String)] = MapPartitions
RDD[29] at mapValues at <console>:46
```

Now here is a new implementation that uses *pattern matching*.

I've also made two other additions, the solutions to the last exercises, which remove empty words "" and fix mixed capitalization, using the following additions:

- `filter(word => word.size > 0)` to remove the empty words. (In Spark and Scala collections, `filter` has the positive sense; what should be retained?) It's indicated by the comment `// #1`.
- `word.toLowerCase` to convert all words to lower case uniformly, so that words like HAMLET, Hamlet, and hamlet in the original texts are treated as the same, since we're counting word occurrences. See comment `// #2`.

```
In [38]: val i11 = sc.wholeTextFiles(shakespeare.toString).
  flatMap {
    case (location, contents) =>
      val words = contents.split(" "\W+"").
        filter(word => word.size > 0) // #1
      val fileName = location.split(pathSeparator).last
      words.map(word => ((word.toLowerCase, fileName), 1)) // #2
  }.
  reduceByKey((count1, count2) => count1 + count2).
  map {
    case ((word, fileName), count) => (word, (fileName, count))
  }.
  groupByKey.
  sortByKey(ascending = true).
  mapValues { iterable =>
    val vect = iterable.toVector.sortBy {
      case (fileName, count) => (-count, fileName)
    }
    vect.mkString(",")
  }
```

```
i11: org.apache.spark.rdd.RDD[(String, String)] = MapPartitionsRDD[39] at
mapValues at <console>:48
```

```
Out[38]: i11: org.apache.spark.rdd.RDD[(String, String)] = MapPartitionsRDD[39] at
mapValues at <console>:48
```

Compare with your exercise solutions above. I added the filtering inside the function passed to `flatMap`. My choice reduces the number of output records from `flatMap` by at most one record per input line, which shouldn't have a significant impact on performance. Filtering itself adds some extra overhead.

Also, the way Spark implements steps like `map`, `flatMap`, `filter`, it would incur about the same overhead if I used `RDD.filter` instead. Note that we could also do the filtering later in the pipeline, after `groupByKey`, for example. So, whichever approach you implemented above is probably fine. You could do performance profiling of the different approaches, but you may not notice a significance difference until you use very large input data sets.

Let's verify we still get reasonable results. Now I'll use Spark's [DataFrame](http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.sql.DataFrame) (<http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.sql.DataFrame>) API for its convenient display options. `DataFrames` are part of [Spark SQL](http://spark.apache.org/docs/latest/sql-programming-guide.html) (<http://spark.apache.org/docs/latest/sql-programming-guide.html>).

First, we need to create an instance of [SQLContext](http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.sql.SQLContext) (<http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.sql.SQLContext>) that we need to access these features.

```
In [39]: import org.apache.spark.sql.SQLContext
```

```
In [40]: val sqlContext = new SQLContext(sc)
```

```
sqlContext: org.apache.spark.sql.SQLContext = org.apache.spark.sql.SQLContext@7ccd9295
```

```
warning: there was one deprecation warning; re-run with -deprecation for details
```

```
Out[40]: sqlContext: org.apache.spark.sql.SQLContext = org.apache.spark.sql.SQLContext@7ccd9295
```

Now, we convert the RDD to a DataFrame with `sqlContext.createDataFrame`, then use `toDF` (convert to another DataFrame ?) with new names for each "column".

```
In [41]: val i1DF = sqlContext.createDataFrame(i1).toDF("word", "locations_counts")
```

```
i1DF: org.apache.spark.sql.DataFrame = [word: string, locations_counts: string]
```

```
Out[41]: i1DF: org.apache.spark.sql.DataFrame = [word: string, locations_counts: string]
```

The `%%dataframe` cell magic provides a nice table layout display.

```
In [42]: %%dataframe
         i1DF
```

```
Out[42]:
```

	word	locations_counts
	a	(loveslabourslost,507),(merrywivesofwindsor,494),(muchadoaboutnothing,492),(asyoulikeit,461),(tamingoftheshrew,445),(twelfthnight,416),(midsummersnightsdream,281),(comedyoferrors,254)
	abandon	(asyoulikeit,4),(tamingoftheshrew,1),(twelfthnight,1)
	abate	(loveslabourslost,1),(midsummersnightsdream,1),(tamingoftheshrew,1)
	abatement	(twelfthnight,1)
	abbess	(comedyoferrors,8)
	abbey	(comedyoferrors,9)
	abominable	(loveslabourslost,1)
	abbreviated	(loveslabourslost,1)
	abed	(asyoulikeit,1),(twelfthnight,1)
	abetting	(comedyoferrors,1)

Okay, now let's explore the new implementation. I start off as before, by calling `wholeTextFiles` :

```
val ii = sc.wholeTextFiles(shakespeare.toString).
```

The function I pass to `flatMap` now looks like this:

```
flatMap {
  case (location, contents) =>
    val words = contents.split(" "\\W+").
      filter(word => word.size > 0) // #
1
    val fileName = location.split(pathSeparator).last
    words.map(word => ((word.toLowerCase, fileName), 1)) // #
2
}.

```

Compare it to the previous version (ignoring the enhancements for blank words and capitalization, marked with the #1 and #2 comments):

```
flatMap { location_contents_tuple2 =>
  val words = location_contents_tuple2._2.split(" "\\W+")
  val fileName = location_contents_tuple2._1.split(pathSeparator)
    .last
  words.map(word => ((word, fileName), 1))
}.

```

Instead of `location_contents_tuple2` a variable name for the whole tuple, I wrote `case (location, contents)`. The `case` keyword says I want to *pattern match* on the object passed to the function. If it's a two-element tuple (and I know it always will be in this case), then *extract* the first element and assign it to a variable named `location` and extract the second element and assign it to a variable named `contents`.

Now, instead of accessing the location and content with the slightly obscure and verbose `location_contents_tuple2._1` and `location_contents_tuple2._2`, respectively, I use meaningful names, `location` and `contents`. The code becomes more concise and more readable.

I'll explore more pattern matching features below.

The `reduceByKey` step is unchanged:

```
reduceByKey((count1, count2) => count1 + count2).
```

To be clear, this isn't a pattern-matching expression; there is no `case` keyword. It's just a "regular" function that takes two arguments, for the two things I'm adding.

My favorite improvement is the next line:

```
map {
  case ((word, fileName), count) => (word, (fileName, count))
}.

```

Compare it to the previous, obscure version:

```
map { word_file_count_tup3 =>
  (word_file_count_tup3._1._1, (word_file_count_tup3._1._2, word_
file_count_tup3._2))
}.
```

The new implementation makes it clear what I'm doing; just shifting parentheses! That's all it takes to go from the `(word, fileName)` keys with `count` values to `word` keys and `(fileName, count)` values. Note that pattern matching works just fine with nested structures, like `((word, fileName), count)`.

I hope you can appreciate how elegant and concise this expression is! Note how I thought of the next transformation I needed to do in preparation for the final group-by, to switch from `((word, fileName), count)` to `(word, (fileName, count))` and *I just wrote it down exactly as I pictured it!*

Code like this makes writing Scala Spark code a sublime experience for me. I hope it will for you, too ;)

The next two expressions are unchanged:

```
groupByKey.
sortByKey(ascending = true).
```

The final `mapValues` now uses pattern matching to sort the `Vector` in each record:

```
mapValues { iterable =>
  val vect = iterable.toVector.sortBy {
    case (fileName, count) => (-count, fileName)
  }
  vect.mkString(", ")
}
```

Compared to the original version, it's again easier to read:

```
mapValues { iterable =>
  val vect = iterable.toVector.sortBy { file_count_tup2 =>
    (-file_count_tup2._2, file_count_tup2._1)
  }
  vect.mkString(", ")
}
```

The function I pass to `sortBy` returns a tuple used for sorting, with `-count` to force *descending* numerical sort (biggest first) and `fileName` to secondarily sort by the file name, for equivalent counts. I could ignore file name order and just return `-count` (not a tuple). However, if you need more repeatable output in a distributed system like Spark, say for example to use in unit test validation, then the secondary sorting by file name is handy.

## Our Final Version: Supporting SQL Queries

To play with some more Spark, let's write SQL queries to explore the resulting data.

To do this, let's first refine the output. Instead of creating a string for the list of `(location, count)` pairs, which is opaque to our SQL schema (i.e., just a string), let's "unzip" the collection into two arrays, one for the `locations` and one for the `counts`. That way, if we ask for the first element of each array, we'll have nicely separate fields that work better with Spark SQL queries.

"Zipping" and "unzipping" work like a mechanical zipper. If I have a collection of tuples, say `List[(String, Int)]`, I convert this single collection of "zippered" values into two collections (in a tuple) of single values, `(List[String], List[Int])`. Zipping is the inverse operation.

Here is our final implementation, `ii1` rewritten with this change.

```
In [43]: val ii = sc.wholeTextFiles(shakespeare.toString).
  flatMap {
    case (location, contents) =>
      val words = contents.split(" "\\W+").
        filter(word => word.size > 0) // #1
      val fileName = location.split(pathSeparator).last
      words.map(word => ((word.toLowerCase, fileName), 1)) // #2
  }.
  reduceByKey((count1, count2) => count1 + count2).
  map {
    case ((word, fileName), count) => (word, (fileName, count))
  }.
  groupByKey.
  sortByKey(ascending = true).
  map {
    case (word, iterable) => // Must use map now, because we'll format
      // Hence, pattern match on the whole input record.

      val vect = iterable.toVector.sortBy {
        case (fileName, count) => (-count, fileName)
      }

      // Use `Vector.unzip`, which returns a single, two element tuple, where
      // element is a collection, one for the locations and one for the counts.
      // I use pattern matching to extract these two collections into variables.
      val (locations, counts) = vect.unzip

      // Lastly, I'll compute the total count across all locations and return
      // a new record with all four fields. The `reduceLeft` method takes
      // that knows how to "reduce" the collection down to a final value,
      // from the left.
      val totalCount = counts.reduceLeft((n1,n2) => n1+n2)

      (word, totalCount, locations, counts)
  }
```

```
ii: org.apache.spark.rdd.RDD[(String, Int, scala.collection.immutable.Vector[String], scala.collection.immutable.Vector[Int])] = MapPartitionsRDD[55] at map at <console>:49
```

```
Out[43]: ii: org.apache.spark.rdd.RDD[(String, Int, scala.collection.immutable.Vector[String], scala.collection.immutable.Vector[Int])] = MapPartitionsRDD[55] at map at <console>:49
```

We've changed the ending `mapValues` call to a `map` call, because we'll construct entirely new records, not just new values with the same keys. Hence the full records, two-element tuples are passed in, rather than just the values, so we'll pattern match on the tuple:

```
map {
  // Must use map now, because we'll
  // format new records.
  case (word, iterable) => // Hence, pattern match on the whole
    // input record.

    val vect = iterable.toVector.sortBy {
      case (fileName, count) => (-count, fileName)
    }
}
```

We have a `Vector[String, Int]` of two-element tuples `(fileName, count)`. We use `Vector.unzip` to create a single, two element tuple, where each element is now a collection, one for the locations and one for the counts. The type is `(Vector[String], Vector[Int])`.

We can also use pattern matching with assignment! We immediately decompose the two-element tuple:

```
// I use pattern matching to extract these two collections
into variables.
val (locations, counts) = vect.unzip
```

Finally, it's convenient to know how many locations and counts we have, so we'll compute another new column for the their count and format a four-element tuple as the final output.

```
// Lastly, I'll compute the total count across all location
s and return
// a new record with all four fields. The `reduceLeft` meth
od takes a function
// that knows how to "reduce" the collection down to a fina
l value, working
// from the left.
val totalCount = counts.reduceLeft((n1,n2) => n1+n2)

(word, totalCount, locations, counts)
}
```

Okay! Now let's create a [DataFrame](http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.sql.DataFrame)

(<http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.sql.DataFrame>) with this data. The `toDF` method just returns the same `DataFrame`, but with appropriate names for the columns, instead of the synthesized names that `createDataFrame` generates (e.g., `_c1`, `_c2`, etc.)

Caching the `DataFrame` in memory prevents Spark from recomputing `ii` from the input files every time I write a query!

Finally, to use SQL, I need to "register" a temporary table.

```
In [44]: val iiDF = sqlContext.createDataFrame(ii).toDF("word", "total_count", "loca
iiDF.cache
iiDF.registerTempTable("inverted_index")
```

```
iiDF: org.apache.spark.sql.DataFrame = [word: string, total_count: int
... 2 more fields]
```

```
warning: there was one deprecation warning; re-run with -deprecation for
details
```

```
Out[44]: iiDF: org.apache.spark.sql.DataFrame = [word: string, total_count: int
... 2 more fields]
```

Let's remind ourselves of the schema:



```
In [45]: iidF.printSchema
```

```
root
|-- word: string (nullable = true)
|-- total_count: integer (nullable = false)
|-- locations: array (nullable = true)
|   |-- element: string (containsNull = true)
|-- counts: array (nullable = true)
|   |-- element: integer (containsNull = false)
```

The following SQL query extracts the top location by count for each word, as well as the total count across all locations for the word. The Spark SQL dialect supports Hive SQL syntax for extracting elements from arrays, maps, and structs ([details](https://cwiki.apache.org/confluence/display/Hive/LanguageManual+UDF#LanguageManualUDF-CollectionFunctions) (<https://cwiki.apache.org/confluence/display/Hive/LanguageManual+UDF#LanguageManualUDF-CollectionFunctions>)). Here I access the first element (index zero) from each array.

```
In [46]: %%SQL
SELECT word, total_count, locations[0] AS top_location, counts[0] AS top_co
FROM inverted_index
```

```
+-----+-----+--...
```

```
Out[46]: +-----+-----+-----+-----+
|      word|total_count|  top_location|top_count|
+-----+-----+-----+-----+
|      a|      3350|loveslabourslost|      507|
|  abandon|         6|  asyoulikeit|         4|
|    abate|         3|loveslabourslost|         1|
| abatement|         1|  twelfthnight|         1|
|    abbess|         8|  comedyoferrors|         8|
|    abbey|         9|  comedyoferrors|         9|
|abbominable|         1|loveslabourslost|         1|
|abbreviated|         1|loveslabourslost|         1|
|    abed|         2|  asyoulikeit|         1|
|  abetting|         1|  comedyoferrors|         1|
+-----+-----+-----+-----+
only showing top 10 rows
```

Unfortunately, the output formatting for the `%%SQL` "cell magic" is not configurable. The `%%DataFrame` magic handles variable width layout and also provides more display options. First, to see its options:

```
In [47]: %%dataframe
```

```
Out[47]: %%dataframe [arguments]
DATAFRAME_CODE
```

DATAFRAME\_CODE can be any numbered lines of code, as long as the last line is a reference to a variable which is a DataFrame.

Option	Description
-----	-----
--limit	The number of records to return (default: 10)
--output	The type of the output: html, csv, json (default: html)

Now here's the previous query again, with the a `WHERE` clause added for good measure:

```
In [48]: val topLocations = sqlContext.sql("""
        SELECT word, total_count, locations[0] AS top_location, counts[0] AS t
        FROM inverted_index
        WHERE word LIKE '%love%' OR word LIKE '%hate%'
        """)
```

```
topLocations: org.apache.spark.sql.DataFrame = [word: string, total_coun
t: int ... 2 more fields]
```

```
Out[48]: topLocations: org.apache.spark.sql.DataFrame = [word: string, total_coun
t: int ... 2 more fields]
```

Now use the `%%dataframe magic`.

```
In [49]: %%dataframe --limit 100
topLocations
```

```
Out[49]:
```

word	total_count	top_location	top_count
beloved	11	tamingoftheshrew	4
cloven	1	loveslabourslost	1
cloves	1	loveslabourslost	1
glove	3	loveslabourslost	2
glover	1	merrywivesofwindsor	1
gloves	5	merrywivesofwindsor	3
hate	22	midsummersnightsdream	9
hated	6	midsummersnightsdream	4
hateful	5	midsummersnightsdream	3
hates	5	asyoulikeit	2
hateth	1	midsummersnightsdream	1
love	662	loveslabourslost	121
loved	38	asyoulikeit	13
lovely	15	midsummersnightsdream	7
lover	33	asyoulikeit	14
lovers	31	midsummersnightsdream	17
loves	51	muchadoaboutnothing	10
lovest	8	tamingoftheshrew	3
loveth	2	loveslabourslost	1
unloved	1	midsummersnightsdream	1
whate	4	tamingoftheshrew	3
whatever	1	tamingoftheshrew	1

A *natural language processing* (NLP) expert might tell you that *love*\_, *\_loved*\_, *\_loves*\_, etc. are really the same word, because they are different conjugations of the verb *\_to love* and *love* is a noun, too. Similarly, should *gloves* (plural) and *glove* (singular) be handled differently?

What we really should do is extract the *stems* of these words and use those instead. NLP toolkits handle this *stemming* for you.

There's also a useful `show` method on `DataFrames`.

```
In [50]: topLocations.show
```

```
+-----+-----+-----+-----+
|  word|total_count|  top_location|top_count|
+-----+-----+-----+-----+
|beloved|      11|tamingoftheshrew|      4|
|cloven|       1|loveslabourslost|      1|
|cloves|       1|loveslabourslost|      1|
|glove|        3|loveslabourslost|      2|
|glover|        1|merrywivesofwindsor|      1|
|gloves|        5|merrywivesofwindsor|      3|
|hate|      22|midsummersnightsd...|      9|
|hated|        6|midsummersnightsd...|      4|
|hateful|       5|midsummersnightsd...|      3|
|hates|        5|asyoulikeit|      2|
|hateth|        1|midsummersnightsd...|      1|
|love|     662|loveslabourslost|     121|
|loved|       38|asyoulikeit|      13|
|lovely|       15|midsummersnightsd...|       7|
|lover|       33|asyoulikeit|      14|
|lovers|       31|midsummersnightsd...|      17|
|loves|       51|muchadoaboutnothing|      10|
|lovest|        8|tamingoftheshrew|       3|
|loveth|        2|loveslabourslost|       1|
|unloved|        1|midsummersnightsd...|       1|
+-----+-----+-----+-----+
```

only showing top 20 rows

By default, it truncates column widths and only prints 20 rows. You can override both:

```
In [51]: topLocations.show(numRows = 40, truncate = false)
```

word	total_count	top_location	top_count
beloved	11	tamingoftheshrew	4
cloven	1	loveslabourslost	1
cloves	1	loveslabourslost	1
glove	3	loveslabourslost	2
glover	1	merrywivesofwindsor	1
gloves	5	merrywivesofwindsor	3
hate	22	midsummersnightsdream	9
hated	6	midsummersnightsdream	4
hateful	5	midsummersnightsdream	3
hates	5	asyoulikeit	2
hateth	1	midsummersnightsdream	1
love	662	loveslabourslost	121
loved	38	asyoulikeit	13
lovely	15	midsummersnightsdream	7
lover	33	asyoulikeit	14
lovers	31	midsummersnightsdream	17
loves	51	muchadoaboutnothing	10
lovest	8	tamingoftheshrew	3
loveth	2	loveslabourslost	1
unloved	1	midsummersnightsdream	1
whate	4	tamingoftheshrew	3
whatever	1	tamingoftheshrew	1

### Note: Named Parameters

I used *named parameters* here, `show(numRows = 40, truncate = false)`, for legibility. They are optional in Scala, as long as you pass the values in the same order as the parameters are declared. You can also use named parameters to write the arguments in any order you want, not just declaration order. So, I could have just written `(40, false)`, but then you would rightly wonder what `false` means in this context.

### Exercises:

See the [Appendix](#) for the solutions to the first two exercises.

- The `glove`, `gloves`, `whate` and `whatever` aren't really the `love` and `hate` we wanted ;) How might you change the query so be more specific.
- Modify the query to return the top two locations and counts.
- Before moving on, try writing other queries. Edit the query in the following cell:

```
In [52]: val sql1 = sqlContext.sql("""
        SELECT * FROM inverted_index
        """)
        sql1.show(10, false)
```

```
+-----+-----+-----+
+-----+-----+-----+
+-----+-----+-----+
|word      |total_count|locations
|counts                                         |
+-----+-----+-----+
+-----+-----+-----+
|a          |3350       |[loveslabourslost, merrywivesofwindsor, muchadoa
boutnothing, asyoulikeit, tamingoftheshrew, twelfthnight, midsummersnight
sdream, comedyoferrors]|[507, 494, 492, 461, 445, 416, 281, 254]|
|abandon    |6          |[asyoulikeit, tamingoftheshrew, twelfthnight]
|[4, 1, 1]                                     |
|abate      |3          |[loveslabourslost, midsummersnightsdream, taming
oftheshrew]
|[1, 1, 1]                                     |
|abatement  |1          |[twelfthnight]
|[1]                                              |
|abbess     |8          |[comedyoferrors]
|[8]                                              |
|abbey      |9          |[comedyoferrors]
|[9]                                              |
|abbominable|1          |[loveslabourslost]
|[1]                                              |
|abbreviated|1          |[loveslabourslost]
|[1]                                              |
|abed       |2          |[asyoulikeit, twelfthnight]
|[1, 1]                                         |
|abetting   |1          |[comedyoferrors]
|[1]                                         |
+-----+-----+-----+
+-----+-----+-----+
+-----+-----+-----+
only showing top 10 rows
```

```
sql1: org.apache.spark.sql.DataFrame = [word: string, total_count: int
... 2 more fields]
```

```
Out[52]: sql1: org.apache.spark.sql.DataFrame = [word: string, total_count: int
... 2 more fields]
```

### Removing the "Stop Words"

Did you notice that one record we saw above was for the word "a". Not very useful if you're using this data for text searching, `_sentiment mining_`, etc. So called `_stop words_`, like `a`, `_an_`, `_the_`, `_he_`, `_she_`, `_it_`, etc., could also be removed.

Recall the `filter` logic I added to remove "", `word => word.size > 0`. I could replace it with `word => keep(word)`, where `keep` is a method that does any additional filtering I want, like removing stop words.

### Exercise:

- Implement the `keep(word: String): Boolean` method and change the `filter` function to use it. Have `keep` return `false` for a small, hard-coded list of stop words (make up your own list or search for one). (See the [Appendix](#) for the solution.)

## More on Pattern Matching Syntax

We've only scratched the surface of pattern matching. Let's explore it some more.

Here's another anonymous function using pattern matching that extends the previous function we passed to `flatMap`:

```
{
  case (location, "") =>
    Array.empty[(String, String), Int] // Return an empty array
  case (location, contents) =>
    val words = contents.split(" "\\W+"")
    val fileName = location.split(pathSep).last
    words.map(word => (word, fileName), 1)
}
```

You can have multiple `case` clauses, some of which might match on specific literal values (" in this case) and others which are more general. The first case clause handles files with no content. The second clause is the same as before.

Pattern matching is *eager*. The first successful match in the order as written will win. If you reversed the order here, the `case (location, "")` would never match and the compiler would throw an "unreachable code" warning for it.

Note that you don't have to put the lines after the `=>` inside braces, `{...}` (although you can). The `=>` and `case` keywords (or the final `}`) are sufficient to mark these blocks. Also, for a single-expression block, like the one for the first case clause, you can put the expression on the same line after the `=>` if you want (and it fits).

Finally, if none of the case clauses matches, then a [MatchError](http://www.scala-lang.org/api/current/index.html#scala.MatchError) (<http://www.scala-lang.org/api/current/index.html#scala.MatchError>) exception is thrown. In our case, we *always* know we'll have two-element tuples, so the examples so far are fine.

Here's a final contrived example to illustrate what's possible, using a sequence of objects of different types:

```
In [53]: val stuff = Seq(1, 3.14159, 2L, 4.4F, ("one", 1), (404F, "boo"), ((11, 12),
stuff.foreach {
  case i: Int           => println(s"Found an Int:    $i")
  case l: Long          => println(s"Found a Long:   $l")
  case f: Float         => println(s"Found a Float:  $f")
  case d: Double        => println(s"Found a Double: $d")
  case (x1, x2) =>
    println(s"Found a two-element tuple with elements of arbitrary type")
  case ((x1a, x1b), _, x3) =>
    println(s"Found a three-element tuple with 1st and 3th elements: ($
  case default         => println(s"Found something else: $default")
}
```

```
Found an Int:    1
Found a Double: 3.14159
Found a Long:   2
Found a Float:  4.4
Found a two-element tuple with elements of arbitrary type: (one, 1)
Found a two-element tuple with elements of arbitrary type: (404.0, boo)
Found a three-element tuple with 1st and 3th elements: (11, 12) and 31
Found something else: hello
```

```
stuff: Seq[Any] = List(1, 3.14159, 2, 4.4, (one,1), (404.0,boo), ((11,1
2),21,31), hello)
```

```
Out[53]: stuff: Seq[Any] = List(1, 3.14159, 2, 4.4, (one,1), (404.0,boo), ((11,1
2),21,31), hello)
```

A few notes.

- A literal like `1` is inferred to be `Int`, while `3.14159` is inferred to be `Double`. Add `L` or `F`, respectively, to infer `Long` or `Float` instead.
- Note how we mixed specific type checking, e.g., `i: Int`, with more loosely-typed expressions, e.g., `(x1, x2)`, which expects a two-element tuple, but the element types are unconstrained.
- All the words `i`, `l`, `f`, `d`, `x1`, `x2`, `x3`, and `default` are arbitrary variable names. Yes `default` is not a keyword, but an arbitrary choice for a variable name. We could use anything we want.
- The last `default` clause specifies a variable with no type information. Hence, it matches `_anything_`, which is why this clause must appear last. This is the idiom to use when you aren't sure about the types of things you're matching against and you want to avoid a possible [MatchError](http://www.scala-lang.org/api/current/index.html#scala.MatchError) (<http://www.scala-lang.org/api/current/index.html#scala.MatchError>).
- If you want to match that something `_exists_`, but you don't need to bind it to a variable, then use `_`, as in the three-element tuple example.
- The three-element tuple example also demonstrates that arbitrary nesting of expressions is supported, where the first element is expected to be a two-element tuple.

All the anonymous functions we've seen that use these pattern matching clauses have this format:



```
{
  case firstCase => ...
  case secondCase => ...
  ...
}
```

This format has a special name. It's called a *partial function*. All that means is that we only "promise" to accept arguments that match at least one of our `case` clauses, not any possible input.

The other kind of anonymous function we've seen is a `_total function_`, to be precise.

Recall we said that for total functions you can use either `(...)` or `{...}` around them, depending on the "look" you want. For *partial functions*, you *must* use `{...}`.

Also, recall that we used pattern matching with assignment:

```
val (locations, counts) = vect.unzip
```

`Vector.unzip` (<http://www.scala-lang.org/api/current/#scala.collection.immutable.Vector>) returns a two-element tuple, where each element is a collection. We matched on that tuple and assigned each piece to a variable. Here's another contrived example, with nested tuple elements:

```
In [54]: val (a, (b, (c1, c2), d)) = ("A", ("B", ("C1", "C2"), "D"))
println(s" $a, $b, $c1, $c2, $d")
```

```
A, B, C1, C2, D
```

```
a: String = A
b: String = B
c1: String = C1
c2: String = C2
d: String = D
```

```
Out[54]: d: String = D
```

Try adding an "E" element to the tuple on the right-hand side, without changing the left-hand side. What happens? Try removing the "D" and "E" elements. What happens now?

We'll come back to one last example of pattern matching when we discuss *case classes*.

## Scala's Object Model

Scala is a *hybrid*, object-oriented and functional programming language. The philosophy of Scala is that you exploit object orientation for encapsulation of details, i.e., *modularity*, but use functional programming for its logical precision when implementing those details. Most of what we've seen so far falls into the functional programming camp. Much of data manipulation and analysis is really Mathematics. Functional programming tries to stay close to how functions and values work in Mathematics.

However, when writing non-trivial Spark programs, it's occasionally useful to exploit the object-oriented features.

## Classes vs. Instances

Scala uses the same distinction between classes and instances that you find in Java. Classes are like *templates* used to create instances.

We've talked about the *types* of things, like `word` is a `String` and `totalCount` is an `Int`. A class defines a *type* in the same sense.

Here is an example class that we might use to represent the inverted index records we just created:

```
In [55]: class IIRecord1(
  word: String,
  total_count: Int,
  locations: Array[String],
  counts: Array[Int]) {

  /** CSV formatted string, but use [a,b,c] for the arrays */
  override def toString: String = {
    val locStr = locations.mkString("[", ",", "]") // i.e., "[a,b,c]"
    val cntStr = counts.mkString("[", ",", "]")   // i.e., "[1,2,3]"
    s"$word,$total_count,$locStr,$cntStr"
  }
}

new IIRecord1("hello", 3, Array("one", "two"), Array(1, 2))
```

```
defined class IIRecord1
```

```
Out[55]: IIRecord1 = hello,3,[one,two],[1,2]
```

When defining a class, the argument list after the class name is the argument list for the *primary constructor*. You can define secondary constructors, too, but it's not very common, in part for reasons we'll see shortly.

Note that when you override a method that's defined in a parent class, like Java's `Object.toString`, Scala requires you to add the `override` keyword.

We created an *instance* of `IIRecord1` using `new`, just like in Java.

Finally, as a side note, we've been using `Ints` (integers) all along for the various counts, but really for "big data", we should probably use `Longs`.

## Objects

I've been careful to use the word *instance* for things we create from classes. That's because Scala has built-in support for the [Singleton Design Pattern](https://en.wikipedia.org/wiki/Singleton_pattern) ([https://en.wikipedia.org/wiki/Singleton\\_pattern](https://en.wikipedia.org/wiki/Singleton_pattern)), i.e., when we only want one instance of a class. We use the `object` keyword.

For example, in Java, you define a class with a `static void main(String[] arguments)` method as your entry point into your program. In Scala, you use an `object` to hold `main`, as follows:

```
In [56]: object MySparkJob {  
  
    val greeting = "Hello Spark!"  
  
    def main(arguments: Array[String]) = {  
        println(greeting)  
  
        // Create your SparkContext, etc., etc.  
    }  
}
```

defined object MySparkJob

Just as for classes, the name of the object can be anything you want. There is no `static` keyword in Scala. Instead of adding `static` methods and fields to classes as in Java, you put them in objects instead, as here.

**NOTE:** Because the Scala compiler must generate valid JVM byte code, these definitions are converted into the equivalent, Java-like static definitions in the output byte code.

## Case Classes

Tuples are handy for representing records and for decomposing them with pattern matching. However, it would be nice if the fields were *named*, as well as *typed*. A good use for a class, like our `IIRecord1` above, is to represent this structure and give us named fields. Let's now refine that class definition to exploit some extra, very useful features in Scala.

Consider the following definition of a *case class* that represents our final record type.

```
In [57]: case class IRecord(
  word: String,
  total_count: Int = 0,
  locations: Array[String] = Array.empty,
  counts: Array[Int] = Array.empty) {

  /**
   * Different than our CSV output above, but see toCSV.
   * Array.toString is useless, so format these ourselves.
   */
  override def toString: String =
    s""""IRecord($word, $total_count, $locStr, $cntStr)"""

  /** CSV-formatted string, but use [a,b,c] for the arrays */
  def toCSV: String =
    s"$word,$total_count,$locStr,$cntStr"

  /** Return a JSON-formatted string for the instance. */
  def toJSONString: String =
    s""""{
      |   "word":      "$word",
      |   "total_count": $total_count,
      |   "locations":  ${toJSONArrayString(locations)},
      |   "counts"     ${toArrayString(counts, ", ")}
      | }
    """".stripMargin

  private def locStr = toArrayString(locations)
  private def cntStr = toArrayString(counts)

  // "[" means we don't care what type of elements; we're just
  // calling toString on them!
  private def toArrayString(array: Array[_], delim: String = ","): String =
    array.mkString("[", delim, "]") // i.e., "[a,b,c]"

  private def toJSONArrayString(array: Array[String]): String =
    toArrayString(array.map(quote), ", ")

  private def quote(word: String): String = "\"" + word + "\""
}
```

```
defined class IRecord
```

I said that defining secondary constructors is not very common. In part, it's because I used a convenient feature, the ability to define default values for arguments to methods, including the primary constructor. The default values mean that I can create instances without providing all the arguments explicitly, as long as there is a default value defined, and similarly for calling methods. Consider these two examples:

```
In [58]: val hello = new IIRRecord("hello")
val world = new IIRRecord("world!", 3, Array("one", "two"), Array(1, 2))

println("\n`toString` output:")
println(hello)
println(world)

println("\n`toJSONString` output:")
println(hello.toJSONString)
println(world.toJSONString)

println("\n`toCSV` output:")
println(hello.toCSV)
println(world.toCSV)
```

```
`toString` output:
IIRRecord(hello, 0, [], [])
IIRRecord(world!, 3, [one,two], [1,2])
```

```
`toJSONString` output:
{
  "word":          "hello",
  "total_count": 0,
  "locations":     [],
  "counts"        []
}

{
  "word":          "world!",
  "total_count": 3,
  "locations":     ["one", "two"],
  "counts"        [1, 2]
}
```

```
`toCSV` output:
hello,0,[],[]
world!,3,[one,two],[1,2]
```

```
hello: IIRRecord = IIRRecord(hello, 0, [], [])
world: IIRRecord = IIRRecord(world!, 3, [one,two], [1,2])
```

```
Out[58]: world: IIRRecord = IIRRecord(world!, 3, [one,two], [1,2])
```

I added `toJSONString` to illustrate adding *public* methods, the default visibility, and *private* methods to a class definition. By the way, when there are no methods or non-field variables to define, I can omit the body complete; no empty `{}` required.

Recall that the `override` keyword is required when redefining `toString`.

Okay, what about that `case` keyword? It tells the compiler to do several useful things for us, eliminating a lot of boilerplate that we would have to write for ourselves with other languages, especially Java:

1. Treat each constructor argument as an immutable ( `val` ) private field of the instance.
2. Generate a public reader method for the field with the same name (e.g., `word` ).
3. Generate *correct* implementations of the `equals` and `hashCode` methods, which people often implement incorrectly, as well as a default `toString` method. You can use your own definitions by adding them explicitly to the body. We did this for `toString` , to format the arrays in a nicer way than the default `Array[_].toString` method.
4. Generate an object `IIRRecord` , i.e., with the same name. The object is called the *companion object*.
5. Generate a "factory" method in the companion object that takes the same argument list and instantiates an instance.
6. Generate helper methods in the companion object that support pattern matching.

Points 1 and 2 make each argument behave as if they are public, read-only fields of the instance, but they are actually implemented as described.

Point 3 is important for correct behavior. Case class instances are often used as keys in [Maps](http://www.scala-lang.org/api/current/index.html#scala.collection.Map) (<http://www.scala-lang.org/api/current/index.html#scala.collection.Map>) and [Sets](http://www.scala-lang.org/api/current/index.html#scala.collection.Set) (<http://www.scala-lang.org/api/current/index.html#scala.collection.Set>), Spark RDD and DataFrame methods, etc. In fact, you should *only* use your case classes or Scala's built-in types with well-defined `hashCode` and `equals` methods (like `Int` and other number types, `String` , tuples, etc.) as keys.

For point 4, the *companion object* is generated automatically by the compiler. It adds the "factory" method discussed in point 5, and methods that support pattern matching, point 6. You can explicitly define these methods and others yourself, as well as fields to hold state. The compiler will still insert these other methods. However, see [Ambiguities with Companion Objects](#). The bottom line is that you shouldn't define case classes in notebooks like this with extra methods in the companion object, due to parsing ambiguities.

Point 5 means you actually rarely use `new` when creating instances. That is, the following are effectively equivalent:

```
In [59]: val hello1 = new IIRRecord("hello1")
        val hello2 = IIRRecord("hello2")
```

```
hello1: IIRRecord = IIRRecord(hello1, 0, [], [])
hello2: IIRRecord = IIRRecord(hello2, 0, [], [])
```

```
Out[59]: hello2: IIRRecord = IIRRecord(hello2, 0, [], [])
```

What actually happens in the second case, without `new` ? The "factory" method is actually called `apply` . In Scala, whenever you put an argument list after any `_instance_`, including these objects , as in the `hello2` case, Scala looks for an `apply` method to call. The arguments have to match the argument list for `apply` (number of arguments, types of arguments, accounting for default argument values, etc.). Hence, the `hello2` declaration is really this:

```
In [60]: val hello2b = IIRRecord.apply("hello2b")
```

```
hello2b: IIRRecord = IIRRecord(hello2b, 0, [], [])
```

```
Out[60]: hello2b: IIRRecord = IIRRecord(hello2b, 0, [], [])
```

You can exploit this feature, too, in your other classes. We talked about word stemming above. Suppose you write a stemming library and declare an object for as the entry point. Here, I'll just do something simple; assume a trailing "s" means the word is a plural and remove it (a bad assumption...):

```
In [61]: object stem {
  def apply(word: String): String = word.replaceFirst("s$", "") // insert
}

println(stem("dog"))
println(stem("dogs"))
```

defined object stem

dog  
dog

Note how it looks like I'm calling a function or method named `stem`. Scala allows object and class names to start with a lower case letter.

Finally, point 6 means we can use our custom case classes in pattern matching expressions. I won't go into the methods actually implemented in the companion object and how they support pattern matching. I'll just use the "magic" in the following example that "parses" or previously-defined `hello` and `world` instances.

```
In [62]: Seq(hello, world).map {
  case IIRecord(word, 0, _, _) => s"$word with no occurrences."
  case IIRecord(word, cnt, locs, cnts) =>
    s"$word occurs $cnt times: ${locs.zip(cnts).mkString(", ")}"
}
```

```
Out[62]: Seq[String] = List(hello with no occurrences., world! occurs 3 times: (one,1), (two,2))
```

The first case clause ignores the locations and counts, because I know they will be empty arrays if the total count is 0!

The second case clause uses the `zip` method to put the locations and counts back together. Recall we used `unzip` to create the separate collections.

## Datasets and DataFrames

So far, we've mostly used Spark's RDD API. It's common to use case classes to represent the "schema" of records when working with RDDs, but also with a new type, [Dataset\[T\]](http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.sql.Dataset) (<http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.sql.Dataset>), analogous to `RDD[T]`, where the `T` represents the type of records.

A problem with [DataFrames](http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.sql.DataFrame) (<http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.sql.DataFrame>) is the fact that the fields are untyped until you try to access them. `Datasets` restore the type safety of

RDDs by using a case class as the definition of the schema.

Datasets were introduced in Spark 1.6.0, but they are somewhat incomplete in the 1.6.X releases. In Spark 2.0.0, `Dataset` becomes the "parent" class of `DataFrame`. This means that you'll be encouraged to use the greater type safety of `Dataset`, but you can still use `DataFrame` if you want. Now, `DataFrame` will be the equivalent of `Dataset[Row]`, where [Row](http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.sql.Row) (<http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.sql.Row>) is the loosely-typed representation of the row and its columns.

Let's try it out. But first, we need to import some SparkSQL-related code. Scala lets you import code almost anywhere, whereas Java requires imports at the beginning of source files. Scala also lets you import members of instances, not just the static imports supported by Java.

So, the next cell imports some "implicit" from the [SQLContext](http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.sql.SQLContext) (<http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.sql.SQLContext>) instance already in scope. Unfortunately, due to a scoping ambiguity involving notebooks and the Scala interpreter, we need to assign `sqlContext` to a new variable, *then* import from that:

```
In [63]: val sqlc = sqlContext
import sqlc.implicitly._
```

```
sqlc: org.apache.spark.sql.SQLContext = org.apache.spark.sql.SQLContext@7ccd9295
```

```
Out[63]: sqlc: org.apache.spark.sql.SQLContext = org.apache.spark.sql.SQLContext@7ccd9295
```

We'll explain what "implicit" are [later](#). For now, suffice it to say that they are used to "allow" us to call the `as` method on our `iiDF` `DataFrame`, which converts it to a `Dataset[IIRecord]`.

```
In [64]: val iiDS = iiDF.as[IIRecord]
iiDS
```

```
iiDS: org.apache.spark.sql.Dataset[IIRecord] = [word: string, total_count: int ... 2 more fields]
```

```
Out[64]: org.apache.spark.sql.Dataset[IIRecord] = [word: string, total_count: int ... 2 more fields]
```



```
In [65]: iiDS.show
```

word	total_count	locations	counts
a	3350	[loveslabourslost...	[507, 494, 492, 4...
abandon	6	[asyoulikeit, tam...	[4, 1, 1]
abate	3	[loveslabourslost...	[1, 1, 1]
abatement	1	[twelfthnight]	[1]
abbess	8	[comedyoferrors]	[8]
abbey	9	[comedyoferrors]	[9]
abbominable	1	[loveslabourslost]	[1]
abbreviated	1	[loveslabourslost]	[1]
abed	2	[asyoulikeit, twe...	[1, 1]
abetting	1	[comedyoferrors]	[1]
abhorrible	1	[loveslabourslost]	[1]
abhor	5	[asyoulikeit, com...	[1, 1, 1, 1, 1]
abhors	2	[twelfthnight]	[2]
abide	5	[merrywivesofwind...	[3, 2]
abides	1	[muchadoaboutnoth...	[1]
ability	2	[muchadoaboutnoth...	[1, 1]
abject	2	[comedyoferrors, ...]	[1, 1]
abjure	1	[midsummersnights...	[1]
abjured	2	[tamingoftheshrew...	[1, 1]
able	9	[merrywivesofwind...	[4, 2, 1, 1, 1]

only showing top 20 rows

## "Scala for Spark 102"

We've covered a lot already in this notebook, focusing on the most important topics you need to know about Scala for daily use. Let's call them the "Scala for Spark 101" material.

At this point, I suggest you create a new notebook and play with Spark using what you've learned so far, then come back to this point if you run into something we didn't cover already. Chances are you're ready to learn the next bits of useful Scala, the "102" material.

### Importing Everything in a Package

In Java, `import foo.bar.*;` means import everything in the `bar` package.

In Scala, `*` is actually a legal method name; think of defining multiplication for custom numeric types, like `Matrix`. Hence, this import statement in Scala would be ambiguous. Therefore, Scala uses `_` instead of `*`, `import foo.bar._` (with the semicolon inferred).

Incidentally, what would that `*` method definition look like? Something like this:

```

case class Matrix(rows: Array[Array[Double]]) { // Each row is an
  Array[Double]

  /** Multiply this matrix by another. */
  def *(other: Matrix): Matrix = ...

  /** Add this matrix by another. */
  def +(other: Matrix): Matrix = ...

  ...
}

val row1: Array[Array[Double]] = ...
val row2: Array[Array[Double]] = ...
val m1 = Matrix(rows1)
val m2 = Matrix(rows2)
val m1_times_m2 = m1 * m2
val m1_plus_m2 = m1 + m2

```

## Operator Syntax

Wait!! What's this `m1 * m2` stuff?? Shouldn't it be `m1.*(m2)` . It would be really convenient to use "operator syntax", more precisely called *infix operator notation* for many methods like `*` and `+` here. The Scala parser supports this with a simple relaxation of the rules; when a method takes a single argument, you can omit the period `.` and parentheses `(...)` . Hence the following really is equivalent:

```

val m1_times_m2 = m1.*(m2)
val m1_times_m2 = m1 * m2

```

This convenience can lead to confusing code, especially for beginners to Scala, so use it cautiously.

## Traits

*Traits* are similar to Java 8 `_interfaces_`, used to define abstractions, but with the ability to provide "default" implementations of the methods declared. Unlike Java 8 interfaces, traits can also have fields representing "state" information about instances. There is a blurry line between traits and `_abstract classes_`, again where some member methods or fields are not defined. In both cases, a subtype of a trait and/or an abstract class must define any undefined members if you want to construct instances of it.

So, why have both traits and abstract classes? It's because Java only allows *single inheritance\_*; *there can be only one \_parent* type, which is normally where you would use an abstract class, but Scala lets you "mix in" one or more additional traits (or use a trait as the parent class - yes, confusing). A great example "mix in" trait is one that implements logging. Any "service" type can mix in the logging trait to get "instant" access to this reusable functionality. Schematically, it looks like the following:

```
// Assume severity `Level` and `Logger` types defined elsewhere...
trait Logging {

    def log(level: Level, message: String): Unit = logger.log(level
, message)

    private logger: Logger = ...
}

abstract class Service {
    def run(): Unit    // No body, so abstract!
}

class MyService extends Service with Logging {
    def run(): Unit = {
        log(INFO, "Staring MyService...")
        ...
        log(INFO, "Finished MyService")
    }
}
```

`Unit` is Scala's equivalent to Java's `void`. It actually is a true type with a single return value, unlike `void`, but we use it in the same sense of "nothing useful will be returned".

## Ranges

What if you want some numbers between a start and end value? Use a [Range](http://www.scala-lang.org/api/current/index.html#scala.collection.immutable.Range) (<http://www.scala-lang.org/api/current/index.html#scala.collection.immutable.Range>), which has a nice literal syntax, e.g., `1 until 100`, `2 to 200 by 3`.

The `Range` always includes the lower bound. Using `to` in a `Range` makes it *inclusive* at the upper bound. Using `until` makes it *exclusive* at the upper bound. Use `by` to specify a delta, which defaults to `1`.

```
In [66]: 1 until 10
```

```
Out[66]: scala.collection.immutable.Range = Range(1, 2, 3, 4, 5, 6, 7, 8, 9)
```

```
In [67]: 1 to 10
```

```
Out[67]: scala.collection.immutable.Range.Inclusive = Range(1, 2, 3, 4, 5, 6, 7,
8, 9, 10)
```

```
In [68]: 1 to 10 by 3
```

```
Out[68]: scala.collection.immutable.Range = Range(1, 4, 7, 10)
```

When you need a small test data set to play with Spark, ranges can be convenient.

```
In [69]: val rdd7 = sc.parallelize(1 to 50).
          map(i => (i, i%7)).
          groupBy{ case (i, seven) => seven }.
          sortByKey()
          rdd7.take(7).foreach(println)

(0,CompactBuffer((7,0), (14,0), (21,0), (28,0), (35,0), (42,0), (49,0)))
(1,CompactBuffer((1,1), (8,1), (15,1), (22,1), (29,1), (36,1), (43,1), (50,1)))
(2,CompactBuffer((2,2), (9,2), (16,2), (23,2), (30,2), (37,2), (44,2)))
(3,CompactBuffer((3,3), (10,3), (17,3), (24,3), (31,3), (38,3), (45,3)))
(4,CompactBuffer((4,4), (11,4), (18,4), (25,4), (32,4), (39,4), (46,4)))
(5,CompactBuffer((5,5), (12,5), (19,5), (26,5), (33,5), (40,5), (47,5)))
(6,CompactBuffer((6,6), (13,6), (20,6), (27,6), (34,6), (41,6), (48,6)))

rdd7: org.apache.spark.rdd.RDD[(Int, Iterable[(Int, Int)]] = ShuffledRDD
[103] at sortByKey at <console>:37

Out[69]: rdd7: org.apache.spark.rdd.RDD[(Int, Iterable[(Int, Int)]] = ShuffledRDD
[103] at sortByKey at <console>:37
```

### SparkContext

(<http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.SparkContext>) also has a `range` method that effectively does the same thing as `sc.parallelize(some_range)`.

## Scala Interpreter (REPL) vs. Notebooks vs. Scala Compiler

This notebook has been using a running Scala interpreter, a.k.a. *REPL* ("read, eval, print, loop") to parse the Scala code. The Spark distribution comes with a `spark-shell` script that also lets you use the interpreter from the command line, but without the nice notebook UI.

If you use `spark-shell`, there are a few other behavior changes you should know about.

### Using `:paste` Mode

By default the Scala interpreter treats *each line* you enter separately. This can cause surprises compared to how the Scala *compiler* works, where it treats all the code in the same file in the same context.

For example, the following code, where the expression continues on the second line, is handled successfully by the compiler, but not by the interpreter.

```
(1 to 100)
  .map(i => i*i)
```

the Interpreter thinks it finished parsing the expression when it hit the new line after the literal `Range` (<http://www.scala-lang.org/api/current/index.html#scala.collection.immutable.Range>), `1 to 100`. It then throws an error on the opening `.` on the next line. On the other hand, the compiler keeps compiling, ignoring the new line in this case.

This notebook also does the same thing as the "raw" interpreter, but in some cases, notebooks will use an interpreter command, `:paste` that tells the parser to parse all of the lines that follow together, just like the compiler would parse them, until the "end of input", which you indicate with `CTRL-D`.

You can't experiment with it through this notebook, but your session would look something like this:

```
scala> :paste
// Entering paste mode (ctrl-D to finish)

(1 to 10)
.map(i => i*i)
<CTRL-D>

// Exiting paste mode, now interpreting.

res0: scala.collection.immutable.IndexedSeq[Int] = Vector(1, 4, 9,
16, 25, 36, 49, 64, 81, 100)

scala>
```

## Ambiguities with Companion Objects

As I wrote this notebook, I *wanted* to demonstrate using the companion object `IIRRecord` to define a method explicitly, but this leads to an ambiguity later on in the notebook if you attempt to use this method. The notebook gets confused between the case class and the object.

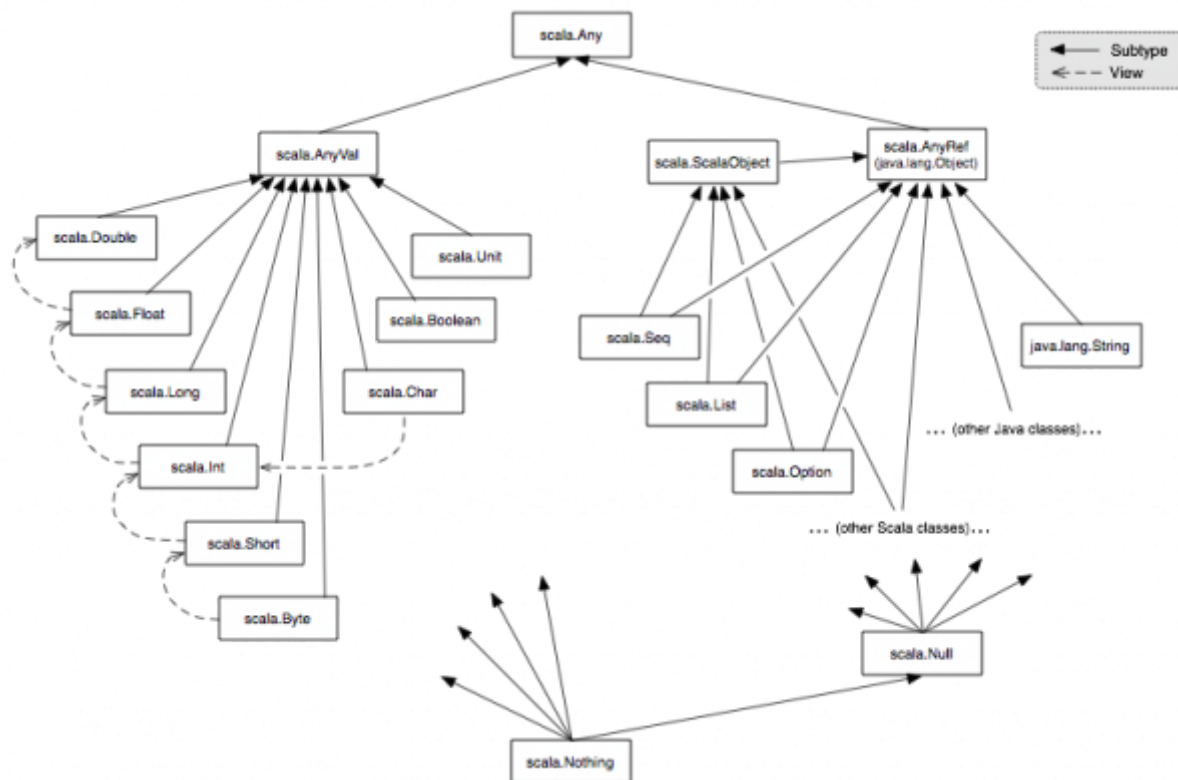
While unfortunate, it's also true that once you start defining more involved case classes, with more than trivial methods and explicit additions to the default companion object, you should really define these types outside the notebook in a compiled library that you use within the notebook.

The details are beyond our scope here, but basically, you set up a project with your Scala code and build it using your favorite build tool. [SBT \(http://www.scala-sbt.org/\)](http://www.scala-sbt.org/) is a popular choice for Scala, but Maven, Gradle, etc. can be used.

You want to generate a `jar` file with the compiled artifacts, then when you start `spark-shell`, submit a Spark job with `spark-submit` or use a notebook environment like this one, you specify the `jar` for inclusion. For `spark-shell` and `spark-submit`, invoke it with the `--jars myproject.jar` option. For Toree with Jupyter, see the discussion on the [FAQ page \(https://toree.incubator.apache.org/documentation/user/faq.html\)](https://toree.incubator.apache.org/documentation/user/faq.html).

## Scala's Type Hierarchy

Scala's type hierarchy is similar to Java's, but with some interesting differences.



In Java, all *reference types* are descended from [java.lang.Object](https://docs.oracle.com/javase/8/docs/api/java/lang/Object.html) (<https://docs.oracle.com/javase/8/docs/api/java/lang/Object.html>). The name *reference type* reflects the fact that the instances for all these types are allocated on the *heap* and program variables are references to those heap locations.

The primitives types, `int`, `long`, etc. are not considered part of the type hierarchy and are treated specially. This is in part a performance optimization, as instances of these types fit in CPU registers and the values are pushed onto stack frames. However, they have wrapper or "boxed" types, `Integer`, `Long`, etc., that are part of the type hierarchy, which you must use with Java's collections, for example (with the exception of arrays).

Instead, Scala treats the primitives at the code level as basically the same as the reference types. You don't use `new Int(100)` for example, but you can call methods on `Int` instances. The code generated, in most cases, uses the optimized JVM primitives.

Hence, the Scala type hierarchy defines a type [Any](http://www.scala-lang.org/api/current/#scala.Any) (<http://www.scala-lang.org/api/current/#scala.Any>) to be the a parent type of *both* reference types and "value" types (for the primitives). Each of those subhierarchies have parent types, [AnyRef](http://www.scala-lang.org/api/current/#scala.AnyRef) (<http://www.scala-lang.org/api/current/#scala.AnyRef>) is effectively the same as [java.lang.Object](https://docs.oracle.com/javase/8/docs/api/java/lang/Object.html) (<https://docs.oracle.com/javase/8/docs/api/java/lang/Object.html>), and [AnyVal](http://www.scala-lang.org/api/current/#scala.AnyVal) (<http://www.scala-lang.org/api/current/#scala.AnyVal>) is the parent of the value types.

Finally, for better "soundness", the Scala type system defines a real type to represent [Null](http://www.scala-lang.org/api/current/#scala.Null) (<http://www.scala-lang.org/api/current/#scala.Null>) and [Nothing](http://www.scala-lang.org/api/current/#scala.Nothing) (<http://www.scala-lang.org/api/current/#scala.Nothing>). By defining `Null` to be the subtype of all reference types `AnyRefs` (but not `AnyVals`), it supports at the type level the (unfortunate) practice of using `null` for a reference value.

However, `null` is not allowed for an `AnyVal`, so the true "bottom type" of the hierarchy is `Nothing`. Why is that useful. I'll explain in the next section.

## Try vs. Option vs. null

Recall the signature of our `curl` method near the beginning of this notebook:

```
def curl(sourceURLString: String, targetDirectoryString: String): File = ...
```

It returns a `File` when everything goes well, but it could throw an exception. An alternative is return a `Try[File]`, where the [Try](http://www.scala-lang.org/api/current/index.html#scala.util.Try) (<http://www.scala-lang.org/api/current/index.html#scala.util.Try>) encapsulates both cases in the return value, as we'll discuss next. We'll also discuss an alternative, [Option](http://www.scala-lang.org/api/current/index.html#scala.Option) (<http://www.scala-lang.org/api/current/index.html#scala.Option>).

Suppose instead that we declared `curl` to return `util.Try[T]` (<http://www.scala-lang.org/api/current/index.html#scala.util.Try>), where `T` is `java.io.File`. The only change to the body would be to simply add `Try` before the opening bracket:

```
def curl(sourceURLString: String, targetDirectoryString: String): Try[File] = Try {...}
```

Now, the reader knows from the method signature that it might fail somehow. If a call fails, the relevant exception will be returned wrapped in a subclass of `Try`, called `util.Failure[T]` (<http://www.scala-lang.org/api/current/index.html#scala.util.Failure>). However, if `curl` succeeds, the `File` will be returned wrapped in the other subclass of `Try`, `util.Success[T]` (<http://www.scala-lang.org/api/current/index.html#scala.util.Success>).

Because of Scala's type safety, the caller of `curl` must determine which result was returned and handle it appropriately. That is, the caller must determine if a `Success` or `Failure` was returned and handle it appropriately.

Scala does not have exception declarations like Java. So, looking at the signature of our original version, there's no obvious way to know if it throws an exception or returns `null` on failure:

```
def curl(sourceURLString: String, targetDirectoryString: String): File = {...}
```

If we choose to catch exceptions internally and return `null`, the caller has to remember to check for `null`. Otherwise, the infamous [NullPointerException](https://docs.oracle.com/javase/8/docs/api/java/lang/NullPointerException.html) (<https://docs.oracle.com/javase/8/docs/api/java/lang/NullPointerException.html>) might happen occasionally if the caller assumes a non-`null` value is returned. So, using `Try[T]` prevents us from this loophole. *It helps the user do the right thing!*

Also, using `Try` rather than simply throwing an exception, means that `curl` always returns "normally", so the caller maintains full control of the call stack and special exception-catching logic isn't required.

What are all the possible valid subclasses of `Try` ? Really, there are only two, `Success` and `Failure` . It would be a mistake to allow a user to define other subtypes, like `MaybeCouldFailButWhoKnows` , because users of `Try` in pattern matching will always want to know that there are only two possibilities. Scala adds a keyword to enforce this logical behavior. `Try` is actually declared as follows:

```
sealed abstract class Try[+T] extends AnyRef
```

(`AnyRef` is the same as Java's `Object` supertype.) The `sealed` keyword says that *no* subclasses of `Try` can be declared, *except* in the same source file (which the library author wrote). Hence, users of `Try` can't declare their own subclasses, subverting the logical structure of this type hierarchy and other user's code that relies on this structure.

What if we have a situation where it makes no sense to involve an exception, but we want the same logically handling? This is where `Option[T]` (<http://www.scala-lang.org/api/current/index.html#scala.Option>) comes in.

`Option` is analogous to `Try` , it is a `sealed abstract` type with two possible subtypes:

- `Some[T]` (<http://www.scala-lang.org/api/current/index.html#scala.None>): I have a an instance of `T` for your, inside the `Some[T]` .
- `None` (<http://www.scala-lang.org/api/current/index.html#scala.None>): I don't have a value for your, sorry.

Note that a hash map is a great example where I either have a value for a given key or I don't. Therefore, for Scala's `Map[K,V]` (<http://www.scala-lang.org/api/current/index.html#scala.collection.Map>) abstraction, where `K` is the key type and `V` is the value type, the `get` method has this signature:

```
def get(key: K): Option[V]
```

One again, you know from the type signature that you may or may not get a value instance for the input key, *and* you **must** determine whether you got a `Some[V]` or a `None` as the result. Once again, we avoid returning a `null` value and risking a `NullPointerException` if we forget to handle it.

So, how do we determine which `Option[T]` was returned? Let's look a few examples using `Option` . Can you guess what they are doing? Check the [Option Scaladocs](http://www.scala-lang.org/api/current/#scala.Option) (<http://www.scala-lang.org/api/current/#scala.Option>) to confirm. `Try` can be used similarly, with a few other ways available that we won't discuss here (but see the [Try Scaladocs](http://www.scala-lang.org/api/current/#scala.util.Try) (<http://www.scala-lang.org/api/current/#scala.util.Try>)).



```
In [70]: val options = Seq(None, Some(2), Some(3), None, Some(5))

options.foreach { o =>
  println(o.getOrElse("None"))
}
```

```
None
2
3
None
5
```

```
options: Seq[Option[Int]] = List(None, Some(2), Some(3), None, Some(5))
```

```
Out[70]: options: Seq[Option[Int]] = List(None, Some(2), Some(3), None, Some(5))
```

```
In [71]: options.foreach {
  case None    => println(None)
  case Some(i) => println(i)  // Note how we extract the enclosed value.
}
```

```
None
2
3
None
5
```

If you just want to ignore the `None` values, use a `_for` comprehension:

```
In [72]: for {
  option <- options  // loop through the options, assign each to "option"
  value  <- option   // extract the value from the Some, or if None, skip
} println(value)
```

```
2
3
5
```

Finally, you might wonder how `None` is declared. Consider this example:

```
In [73]: val opts: Seq[Option[String]] = Seq(Some("hello"), None, Some("world!"))
opts.foreach(println)
```

```
Some(hello)
None
Some(world!)
```

```
opts: Seq[Option[String]] = List(Some(hello), None, Some(world!))
```

```
Out[73]: opts: Seq[Option[String]] = List(Some(hello), None, Some(world!))
```

This works, so it must mean that `None` is a valid subclass of `Option[String]`. That's actually true for all `Option[T]`. How can a single object be a valid subtype for *all* of them? Here is how it's declared (omitting some details):

```
object None extends Option[Nothing] {...}
```

`None` carries no "state" information, because it doesn't wrap an instance like `Some[T]` does. Hence, we only need one instance for all uses, so it's declared as an object. Recall we mentioned above that the type system has a `Nothing` (<http://www.scala-lang.org/api/current/#scala.Nothing>) type, which is a subtype of all other types. Without diving into too many details, if a variable is of type `Option[String]`, then you can use an `Option[Nothing]` for it (i.e., the latter is a subtype of the former). This is why `Nothing` is useful, for cases like `None`, so we can have one instance of it, but still obey the rules of Scala's object-oriented type system.

## Implicits

Scala has a powerful mechanism known as *implicits* that is used in the Spark Scala API. Implicits are a big topic, so we'll focus just on the uses of it that are most important to understand.

### Type Conversions

We used RDD methods like `reduceByKey` above, but if you search for this method in the [RDD Scaladoc page](http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.rdd.RDD) (<http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.rdd.RDD>), you won't find it. Instead it's defined in the [PairRDDFunctions](http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.rdd.PairRDDFunctions) (<http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.rdd.PairRDDFunctions>) type (along with all the other `*ByKey` methods). So, how can we use these methods as if they are defined for `RDD` ??

When the Scala compiler sees code calling a method that doesn't exist on the type, it looks for an *implicit conversion* in the current scope, which can transform the instance into another type (i.e., by wrapping it), where the other type provides the needed method. The full signature inferred for the method as it's used must match the definition in the wrapping class.

**Note:** If you don't find a method in the [Spark Scaladocs](http://spark.apache.org/docs/latest/api/scala/index.html#package) (<http://spark.apache.org/docs/latest/api/scala/index.html#package>) for a type where you think it should be defined, look for related helper types with the method.

Here's a small Scala example of how this works:

```
In [74]: // A sample class. Note it doesn't define a `toJSON` method:
case class Person(name: String, age: Int = 0)
```

```
defined class Person
```

```
In [75]: // To scope them, define implicit conversions within an object
object implicits {

    // `implicit` keyword tells the compiler to consider this conversion.
    // It takes a `Person`, returning a new instance of `PersonToJSONString`
    // then resolves the invocation of `toJSON`.
    implicit class PersonToJSONString(person: Person) {
        def toJSON: String = s""""{"name": ${person.name}, "age": ${person.a
    }

import implicits._           // Now it is visible in the current scope.

val p = Person("Dean Wampler", 39)

// Magic conversion to `PersonToJSONString`, then `toJSON` is called.
p.toJSON
```

```
defined object implicits
p: Person = Person(Dean Wampler,39)
```

```
Out[75]: String = {"name": Dean Wampler, "age": 39}
```

For RDDs, the implicit conversions to `PairRDDFunctions` and other support types are handled for you. However, when you use Spark SQL and the [DataFrame](http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.sql.DataFrame) (<http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.sql.DataFrame>) API, you'll need to import some of these conversions yourself:

```
In [76]: val sqlc = sqlContext
import sqlc.implicits._

sqlc: org.apache.spark.sql.SQLContext = org.apache.spark.sql.SQLContext@7
ccd9295

Out[76]: sqlc: org.apache.spark.sql.SQLContext = org.apache.spark.sql.SQLContext@7
ccd9295
```

```
In [77]: val wtc = iiDF.select($"word", $"total_count")
        wtc.show
```

word	total_count
a	3350
abandon	6
abate	3
abatement	1
abbess	8
abbey	9
abbominable	1
abbreviated	1
abed	2
abetting	1
abhorrible	1
abhor	5
abhors	2
abide	5
abides	1
ability	2
abject	2
abjure	1
abjured	2
able	9

only showing top 20 rows

```
wtc: org.apache.spark.sql.DataFrame = [word: string, total_count: int]
```

```
Out[77]: wtc: org.apache.spark.sql.DataFrame = [word: string, total_count: int]
```

The column-reference syntax  `$"name"`  is implemented using the same mechanism in the Scala library that implements interpolated strings,  `s"$foo"` . The  `import sqlc.implicits._`  makes it available.

Note we imported something from an  `_instance_` , rather than a package or type, as allowed in Java. This can be a useful feature in Scala, but it's also fragile. If you try  `import sqlContext.implicits._` , you'll get a compiler error that a "stable identifier" is required. It turns out that doing the value assignment,  `val sqlc = sqlContext`  first meets this requirement. This is unique to the notebook environment. You normally won't see this problem if you use the  `spark-shell`  that comes with a Spark distribution or you write a Spark program and compile it with the Scala compiler.

However, it would be better if Spark defined this  `implicits`  object on the  `SQLContext`  companion object instead of on instances of it!

For completeness, but unrelated to implicits, the  `DataFrame`  API lets you write SQL-like queries with a programmatic API. If you want to use built in functions like  `min` ,  `max` , etc. on columns, you need the following  `import`  statement:

```
In [78]: import org.apache.spark.sql.functions._
```

Now we can use `min`, `max`, `avg`, etc.

```
In [79]: val mma = iidF.select(min("total_count"), max("total_count"), avg("total_co
mma.show
```

```
+-----+-----+-----+
|min(total_count)|max(total_count)|  avg(total_count)|
+-----+-----+-----+
|              1|          5208|16.651743683350947|
+-----+-----+-----+
```

```
mma: org.apache.spark.sql.DataFrame = [min(total_count): int, max(total_c
ount): int ... 1 more field]
```

```
Out[79]: mma: org.apache.spark.sql.DataFrame = [min(total_count): int, max(total_c
ount): int ... 1 more field]
```

## Implicit Method Arguments

One other use of implicits worth understanding is *implicit arguments* to methods. You will encounter this mechanism used when you read the Spark Scaladocs, even though you might never realize you're actually using it in your code!

Recall I mentioned previously that you can define default values for method arguments. I just used it for the `age` argument for `Person`:

```
case class Person(name: String, age: Int = 0)
```

Sometimes we need something more sophisticated. For example, our library might have a group of methods that need a special argument passed to them that provides useful "context" information, but you don't want the user to be required to explicitly pass this argument every time. Other times you might use implicit arguments to make the API "cleaner", but still have some control over what's allowed.

Here's an example, that's partly inspired by Scala's [Seq.sum](http://www.scala-lang.org/api/current/#scala.collection.Seq) (<http://www.scala-lang.org/api/current/#scala.collection.Seq>) method. Wouldn't it be great if I happen to have a collection of things I can "add" together, if I could just call `sum` on the collection? Let's do this in a slightly different way, with a helper `sum` method outside of `Seq`.

```
In [80]: trait Add[T] {
          def add(t1: T, t2: T): T
        }

        // Nested implicits so they don't conflict with the previous object implicit
        object Adder {
          object implicits {
            implicit val intAdd = new Add[Int] {
              def add(i1: Int, i2: Int): Int = i1+i2
            }
            implicit val doubleAdd = new Add[Double] {
              def add(d1: Double, d2: Double): Double = d1+d2
            }
            implicit val stringAdd = new Add[String] {
              def add(s1: String, s2: String): String = s1+s2
            }
            // etc...
          }
        }

        import Adder.implicits._

        // NOTE: TWO argument lists!
        def sum[T](ts: Seq[T])(implicit adder: Add[T]): T = {
          ts.reduceLeft((t1, t2) => adder.add(t1, t2))
        }
```

```
defined trait Add
defined object Adder
```

```
sum: [T](ts: Seq[T])(implicit adder: Add[T])T
```

```
In [81]: sum(0 to 10)
```

```
Out[81]: Int = 55
```

```
In [82]: sum(0.0 to 5.5 by 0.3)
```

```
Out[82]: Double = 51.299999999999999
```

```
In [83]: sum(Seq("one", "two", "three"))
```

```
Out[83]: String = onetwothree
```

```
In [84]: // Will fail, because there's no Add[Char] in scope:
          sum(Seq('a', 'b', 'c')) // Characters
```

```
Out[84]: Name: Compile Error
          Message: <console>:53: error: could not find implicit value for parameter
          adder: Add[Char]
          sum(Seq('a', 'b', 'c')) // Characters
                ^
```

```
StackTrace:
```

So, the implicit values `intAdd`, `doubleAdd`, and `stringAdd`, were used by the Scala interpreter for the `adder` argument in the second *argument list* for `sum`. Note that you have to use a second argument list and all arguments there must be implicit.

We could have avoided using implicit arguments if we defined custom `sum` methods for every type. That would have been simpler in this trivial case, but for nontrivial methods, the duplication is worth avoiding. Another advantage of this mechanism is that the user can define her own implicit `Add[T]` instances for domain types (say for example, `Money`) and they would "just work".

The Scala collections API uses this mechanism to know how to construct a new collection of the same kind as the input collection when you use `map`, `flatMap`, `reduceLeft`, etc.

Spark uses this pattern for [Encoders](http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.sql.Encoder) (<http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.sql.Encoder>) in Spark SQL. Encoders are used to serialize values into the new, compact memory encoding introduced in the *Tungsten* project (see for example, [here](https://spark-summit.org/2015/events/deep-dive-into-project-tungsten-bringing-spark-closer-to-bare-metal/) (<https://spark-summit.org/2015/events/deep-dive-into-project-tungsten-bringing-spark-closer-to-bare-metal/>)). Here's an example of creating a [Dataset](http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.sql.Dataset) (<http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.sql.Dataset>), where the `toDS` method is first "added" to a Scala [Seq](http://www.scala-lang.org/api/current/#scala.collection.Seq) (<http://www.scala-lang.org/api/current/#scala.collection.Seq>) through an implicit conversion (specifically [SQLImplicits.localSeqToDatasetHolder](http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.sql.SQLImplicits) (<http://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.sql.SQLImplicits>), which is brought into scope by the `import sqlc.implicits._` statement earlier) and then `toDS` uses `Encoders` internally.

```
In [85]: (0 to 10).toDS()
```

```
Out[85]: org.apache.spark.sql.Dataset[Int] = [value: int]
```

## Conclusions

I appreciate the effort you put into studying this notebook. I hope you enjoyed it as much as I enjoyed writing it. Please post issues on how I can improve it to the [GitHub repo](https://github.com/deanwampler/JustEnoughScalaForSpark) (<https://github.com/deanwampler/JustEnoughScalaForSpark>).

Now you know the core elements of Scala that you need for using the Spark Scala API. I hope you can appreciate the power and elegance of Scala. I hope you will choose to use it for all of your data engineering tasks, not just for Spark.

What about data science? There are many people who use Scala for data science in Spark, but today Python and R have much richer libraries for Mathematics and Machine Learning. That will change over time, but for now, you'll need to decide which language best fits your needs.

As you use Scala, there will be more things you'll want to understand that we haven't covered, including common idioms, conventions, and tools used in the Scala community. The references at the beginning of the notebook will give you the information you need.

Best wishes.

[Dean Wampler, Ph.D. \(mailto:deanwampler@gmail.com\)](mailto:deanwampler@gmail.com)  
[@deanwampler \(http://twitter.com/deanwampler\)](http://twitter.com/deanwampler)

## Appendix: Exercise Solutions

Let's discuss the solutions to exercises that weren't already solved earlier in the notebook.

### Filter for Plays that Have "of" in the Name

You can add the condition (comment `// <== here`) immediate after defining `play`. You could do it later, after either of the subsequent two expressions, but then you're doing needless computation. Change `true` to `false` to print plays that don't contain "of".

```
In [86]: val list2 = for {
  play <- plays
  if (play.contains("of") == true)                                // <== here
  playFileString = targetDirName + pathSeparator + play
  playFile = new File(playFileString)
} yield {
  val successString = if (playFile.exists) "Success!" else "NOT FOUND!!"
  "%-40s\t%s".format(playFileString, successString)
}
list2.foreach(println)
```

```
/home/jovyan/work/data/shakespeare/tamingoftheshrew      Success!
/home/jovyan/work/data/shakespeare/comedyoferrors         Success!
/home/jovyan/work/data/shakespeare/merrywivesofwindsor   Success!
```

```
list2: Seq[String] = List(/home/jovyan/work/data/shakespeare/tamingofthes
hrew      Success!, /home/jovyan/work/data/shakespeare/comedyoferrors      S
uccess!, /home/jovyan/work/data/shakespeare/merrywivesofwindsor Success!)
```

```
Out[86]: list2: Seq[String] = List(/home/jovyan/work/data/shakespeare/tamingofthes
hrew      Success!, /home/jovyan/work/data/shakespeare/comedyoferrors      S
uccess!, /home/jovyan/work/data/shakespeare/merrywivesofwindsor Success!)
```

### More Specific "Love" and "Hate" Words

One reasonable choice to prevent seeing `glove`, `whatever`, etc. is to only find words that start with `love` and `have`. Let's also keep `unlove`:



```
In [87]: val topLocationsLoveHate = sqlContext.sql("""
  SELECT word, total_count, locations[0] AS top_location, counts[0] AS t
  FROM inverted_index
  WHERE word LIKE 'love%' OR word LIKE 'unlove%' OR word LIKE 'hate%'
  """)
topLocationsLoveHate.show(40)
```

word	total_count	top_location	top_count
hate	22	midsummersnightsd...	9
hated	6	midsummersnightsd...	4
hateful	5	midsummersnightsd...	3
hates	5	asyoulikeit	2
hateth	1	midsummersnightsd...	1
love	662	loveslabourslost	121
loved	38	asyoulikeit	13
lovely	15	midsummersnightsd...	7
lover	33	asyoulikeit	14
lovers	31	midsummersnightsd...	17
loves	51	muchadoaboutnothing	10
lovest	8	tamingoftheshrew	3
loveth	2	loveslabourslost	1
unloved	1	midsummersnightsd...	1

```
topLocationsLoveHate: org.apache.spark.sql.DataFrame = [word: string, tot
al_count: int ... 2 more fields]
```

```
Out[87]: topLocationsLoveHate: org.apache.spark.sql.DataFrame = [word: string, tot
al_count: int ... 2 more fields]
```

## Return the Top Two Locations and Counts

We used the `DataFrame` API to write a SQL query that returned the top location and count. Adding the next one is straightforward. What do you observe is returned when there isn't a second location and count?

```
In [88]: val topTwoLocations = sqlContext.sql("""
  SELECT word, total_count,
  locations[0] AS first_location, counts[0] AS first_count,
  locations[1] AS second_location, counts[1] AS second_count
  FROM inverted_index
  WHERE word LIKE '%love%' OR word LIKE '%hate%'
  """)
```

```
topTwoLocations: org.apache.spark.sql.DataFrame = [word: string, total_co
unt: int ... 4 more fields]
```

```
Out[88]: topTwoLocations: org.apache.spark.sql.DataFrame = [word: string, total_co
unt: int ... 4 more fields]
```

```
In [89]: topTwoLocations.show(100)
```

```
+-----+-----+-----+-----+-----+
--+-----+
|   word|total_count|   first_location|first_count|   second_locati
on|second_count|
+-----+-----+-----+-----+-----+
--+-----+
| beloved|      11|   tamingoftheshrew|      4|   asyoulike
it|      3|
| cloven|      1|   loveslabourslost|      1|           nu
ll|    null|
| cloves|      1|   loveslabourslost|      1|           nu
ll|    null|
| glove|      3|   loveslabourslost|      2|   twelfthnig
ht|      1|
| glover|      1| merrywivesofwindsor|      1|           nu
ll|    null|
| gloves|      5| merrywivesofwindsor|      3|   asyoulike
it|      1|
|  hate|     22|midsummersnightsd...|      9|   asyoulike
it|      6|
|  hated|      6|midsummersnightsd...|      4|   asyoulike
it|      2|
| hateful|      5|midsummersnightsd...|      3|   loveslabourslo
st|      1|
|  hates|      5|   asyoulikeit|      2| merrywivesofwinds
or|      1|
| hateth|      1|midsummersnightsd...|      1|           nu
ll|    null|
|  love|     662|   loveslabourslost|    121|   asyoulike
it|    119|
|  loved|      38|   asyoulikeit|     13| muchadoaboutnothi
ng|     13|
| lovely|     15|midsummersnightsd...|      7|   tamingoftheshr
ew|      5|
|  lover|     33|   asyoulikeit|     14|midsummersnights
d...|     10|
|  lovers|     31|midsummersnightsd...|     17|           asyoulike
it|      6|
|  loves|     51| muchadoaboutnothing|     10| merrywivesofwinds
or|      9|
| lovest|      8|   tamingoftheshrew|      3| muchadoaboutnothi
ng|      2|
| loveth|      2|   loveslabourslost|      1|   tamingoftheshr
ew|      1|
| unloved|      1|midsummersnightsd...|      1|           nu
ll|    null|
|  whate|      4|   tamingoftheshrew|      3|   asyoulike
it|      1|
|whatever|      1|   tamingoftheshrew|      1|           nu
ll|    null|
+-----+-----+-----+-----+-----+
--+-----+
```

## Removing Stop Words

Recall you were asked to implement a `keep(word: String): Boolean` method that filters stop words.

First, let's implement `keep`. You can find lists of stop words on the web. One such list for English can be found [here]( \* From <http://norm.al/2009/04/14/list-of-english-stop-words/> (<http://norm.al/2009/04/14/list-of-english-stop-words/>)). It includes many words that you might not consider stop words. Nevertheless, I'll just use a smaller list here.

Note that I'll use a Scala [Set](http://www.scala-lang.org/api/current/index.html#scala.collection.immutable.Set) (<http://www.scala-lang.org/api/current/index.html#scala.collection.immutable.Set>) to hold the stop words. We want  $O(1)$  look-up performance. We just want to know if the word is in the set or not.

I'll also add "", so I can remove the explicit test for it.

Finally, we'll embed the whole thing in a new Scala `object`. This extra encapsulation is a way to work around occasional problems with "task not serializable" errors.

**WARNING:** The definition in the next cell may trigger a `Task not serializable` error in the cell that follows, where it is used. Of so, this is "quirk" of the Scala interpreter running with the notebook environment. This code should work without issues in Spark applications that you write, i.e., that you compile into applications with `scalac`.

```

In [90]: object IIStopWords {
    val stopWords = Set("", "a", "an", "and", "I", "he", "she", "it", "the")

    /**
     * If the set contains the word, we return false - we don't want to keep
     * Note we assume the word has already been converted to lower case!
     */
    def keep(word: String): Boolean = stopWords.contains(word) == false

    def compute(sc: org.apache.spark.SparkContext, input: String) = {
      sc.wholeTextFiles(input).
      flatMap {
        case (location, contents) =>
          val words = contents.split("\\W+").
            map(word => word.toLowerCase). // Do this early, before
            filter(word => keep(word)) // <== filter here
          val fileName = location.split(java.io.File.separator).last
          words.map(word => ((word, fileName), 1))
      }.
      reduceByKey((count1, count2) => count1 + count2).
      map {
        case ((word, fileName), count) => (word, (fileName, count))
      }.
      groupByKey.
      sortByKey(ascending = true).
      map {
        case (word, iterable) =>
          val vect = iterable.toVector.sortBy {
            case (fileName, count) => (-count, fileName)
          }
          val (locations, counts) = vect.unzip
          val totalCount = counts.reduceLeft((n1, n2) => n1+n2)
          (word, totalCount, locations, counts)
      }
    }
  }
}

```

defined object IIStopWords

```
In [91]: val iiStopWords = IIStopWords.compute(sc, "/home/jovyan/work/data/shakespea
```

```
Out[91]: Name: org.apache.spark.SparkException
Message: Task not serializable
StackTrace:  at org.apache.spark.util.ClosureCleaner$.ensureSerializable
(ClosureCleaner.scala:403)
    at org.apache.spark.util.ClosureCleaner$.org$apache$spark$util$ClosureC
leaner$$clean(ClosureCleaner.scala:393)
    at org.apache.spark.util.ClosureCleaner$.clean(ClosureCleaner.scala:16
2)
    at org.apache.spark.SparkContext.clean(SparkContext.scala:2326)
    at org.apache.spark.rdd.RDD$$anonfun$flatMap$1.apply(RDD.scala:402)
    at org.apache.spark.rdd.RDD$$anonfun$flatMap$1.apply(RDD.scala:401)
    at org.apache.spark.rdd.RDDOperationScope$.withScope(RDDOperationScope.
scala:151)
    at org.apache.spark.rdd.RDDOperationScope$.withScope(RDDOperationScope.
scala:112)
    at org.apache.spark.rdd.RDD.withScope(RDD.scala:385)
    at org.apache.spark.rdd.RDD.flatMap(RDD.scala:401)
    at IIStopWords$.compute(<console>:65)
    ... 52 elided
Caused by: java.io.NotSerializableException: IIStopWords$
Serialization stack:
    - object not serializable (class: IIStopWords$, value: IIStopWord
s$@a08fa6c)
    - field (class: IIStopWords$$anonfun$3, name: $outer, type: class
IIStopWords$)
    - object (class IIStopWords$$anonfun$3, <function1>)
    at org.apache.spark.serializer.SerializationDebugger$.improveException
(SerializationDebugger.scala:40)
    at org.apache.spark.serializer.JavaSerializationStream.writeObject(Java
Serializer.scala:46)
    at org.apache.spark.serializer.JavaSerializerInstance.serialize(JavaSer
ializer.scala:100)
    at org.apache.spark.util.ClosureCleaner$.ensureSerializable(ClosureClea
ner.scala:400)
```

```
In [92]: iiStopWords.take(100).foreach(println)
```

```
Out[92]: Name: Unknown Error
Message: lastException: Throwable = null
<console>:51: error: not found: value iiStopWords
    iiStopWords.take(100).foreach(println)
    ^
```

StackTrace:

One last thing, we now have `filter(word => keep(word))`, but note how we used `println` in the previous cell to see results. We can do something similar with `filter` and instead write `filter(keep)`.

What does this mean exactly? It tells the compiler "convert the *method* `keep` to a *function* and pass that to `filter`." This works because `keep` already does what `filter` wants, take a single string argument and return a boolean result.

Passing `keep` is actually different than passing `word => keep(word)`, which is an *anonymous* function that *calls* `keep`. We are using `keep` as the function itself, rather than constructing a function that uses `keep`.