

# CS5344 Lab 1

AY2018/2019 Semester 2

The purpose of this lab is to get you started with Spark, and learn how to write, compile, debug and execute a simple Spark program. You will also be tasked to write and submit your own Spark program individually.

1. **Reference programs and documentation** from Spark release are available at  
[https://www.tutorialspoint.com/apache\\_spark/apache\\_spark\\_quick\\_guide.htm](https://www.tutorialspoint.com/apache_spark/apache_spark_quick_guide.htm)  
<https://spark.apache.org/docs/2.2.0/rdd-programming-guide.html>
2. A VirtualBox image of Ubuntu with Spark deployment has been configured for you. **Appendix A** gives the instructions on how to download and install it.
3. Alternatively, you can learn to install a stand-alone Spark-2.2.1 instance on Ubuntu by yourself and set up the environment by following the instructions in **Appendix B**. For **Mac users**, you can refer to  
<https://medium.com/luckspark/installing-spark-2-3-0-on-macos-high-sierra-276a127b8b85>
4. **Appendix C** and **D** are basic guides to help you get started and run your first Spark *WordCount* program in Java and Python respectively.
5. **Task: Write a Spark application (in Java OR Python) to count the number of words that begin with each letter.** In other words, for each letter, count the total number of words that start with that letter. The words need not be unique and you can ignore the letter case and all non-alphabetic characters. Run your program over the input file provided. Upload the following to the Lab1 folder in IVLE.
  - (a) Source program file (with documentation within the code)
  - (b) Output file

## Appendix A. Install VirtualBox and Configure Ubuntu Image


1. Install VirtualBox VM <https://www.virtualbox.org/>

VirtualBox supports Windows, Mac OS, Linux, Solaris.

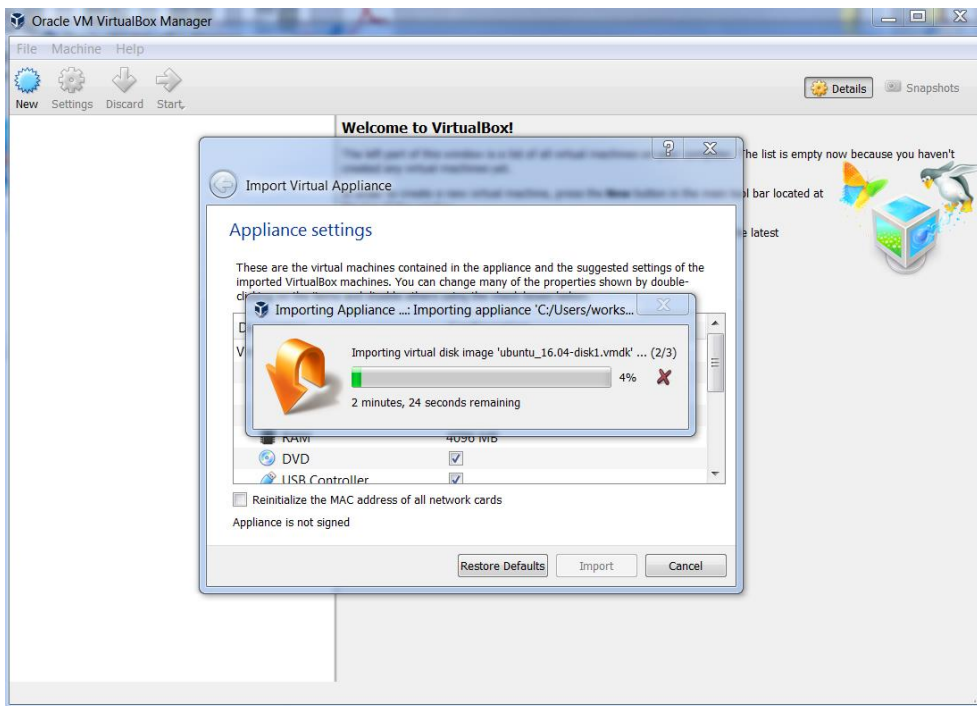
Download the installation file for your operating system. Double-click on file to install it.

2. Download Ubuntu Image from [https://www.dropbox.com/s/i8oopo4upl8mm9x/Ubuntu\\_16.04.zip?dl=0](https://www.dropbox.com/s/i8oopo4upl8mm9x/Ubuntu_16.04.zip?dl=0)

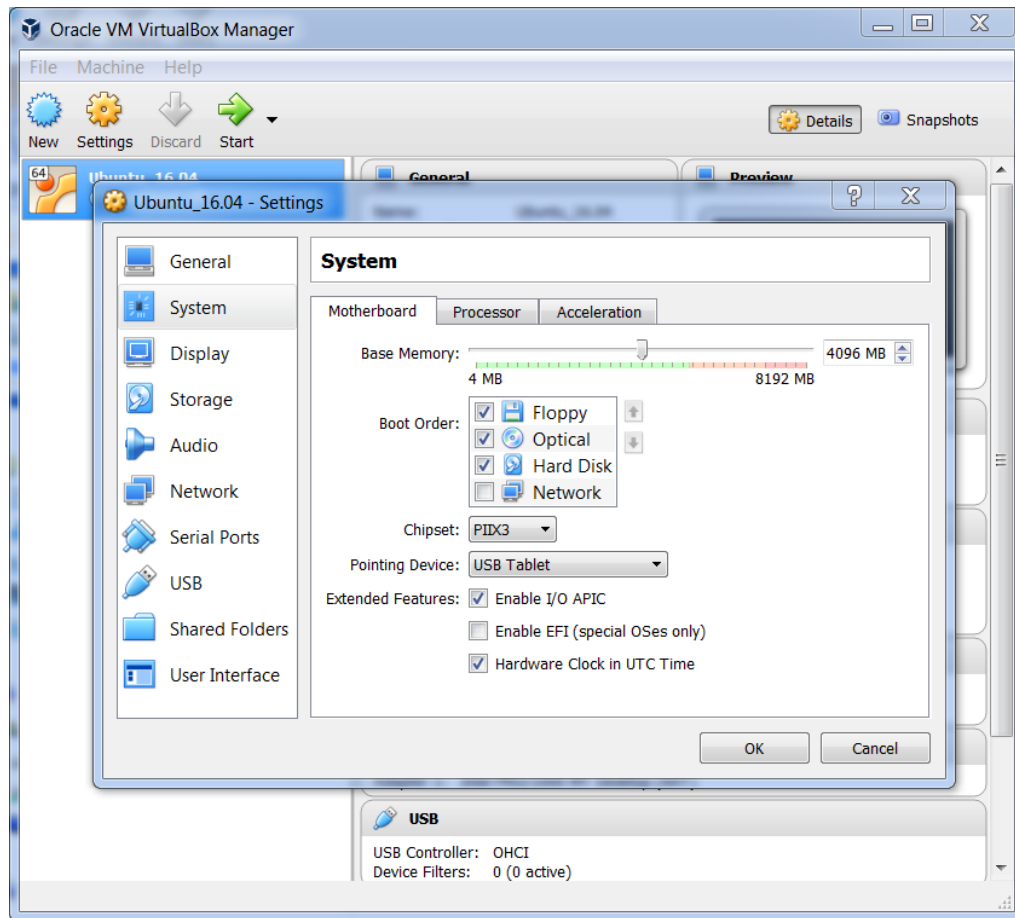
You will see the following file when you unzip ubuntu\_16.04.zip file:

Name	Date modified	Type	Size
 Ubuntu_16.04.ova	15/1/2018 7:41 PM	Open Virtuali...	3,156,076 ...

3. Double click Ubuntu\_16.04.vbox file, and the image will be loaded into VirtualBox.



By default, the VM will take up 4G of your physical memory. It is recommended that the VM memory usage does not exceed half of your total memory. To adjust memory usage, click "Settings" button, click "System" tab, and adjust "Base Memory".



If there is any error while starting the VM related to USB, disable the USB controller from “Settings”, under the “USB” tab.

4. Now you can start the Ubuntu VM. By default, the username is “Spark” and password is “123456” (without quotation marks).

## Appendix B. Install Spark-2.2.1 on Ubuntu-16.04 with JDK 8

### 1. Install JDK 8

Verify Java installation

```
$ java -version
```

If Java is not installed, we install it via the following commands.

```
$ sudo add-apt-repository ppa:webupd8team/java
$ sudo apt-get update && sudo apt-get install oracle-java8-installer
```

It may take some time to download the install. When it is done, set the path as follows.

```
$ sudo gedit /etc/environment
```

Append the following line at the end of the file and save it.

*JAVA\_HOME="/usr/lib/jvm/java-8-oracle"*

### 2. Install Scala

Download Scala in <http://www.scala-lang.org/download/>

#### Other resources

You can find the installer download links for other operating systems, as well as documentation and source code archives for Scala 2.12.4 below.

Archive	System	Size
<a href="#">scala-2.12.4.tgz</a>	Mac OS X, Unix, Cygwin	18.83M
<a href="#">scala-2.12.4.msi</a>	Windows (msi installer)	126.38M
<a href="#">scala-2.12.4.zip</a>	Windows	18.87M
<a href="#">scala-2.12.4.deb</a>	Debian	145.23M
<a href="#">scala-2.12.4.rpm</a>	RPM package	125.81M
<a href="#">scala-docs-2.12.4.tgz</a>	API docs	56.52M
<a href="#">scala-docs-2.12.4.zip</a>	API docs	109.65M
<a href="#">scala-sources-2.12.4.tar.gz</a>	Sources	

```
$ cd /home/Spark/Downloads
$ tar xvf scala-2.12.4.tgz
```

Use the following commands to move the Scala files to the directory /usr/local/scala

```
$ su -
Password:
# cd /home/Spark/Downloads/
# mv scala-2.12.4 /usr/local/scala
# exit
```

If you have not set the password for root account, use the following command to set it.

```
$ sudo passwd
```

Set the path for Scala.

```
$ sudo gedit /etc/environment
```

Append the following clause to the end of PATH = “/usr/local/sbin:.....” in the file, and save it.  
*:/usr/local/scala/bin*

### 3. Install Maven (to compile java files)

Download Maven from <https://maven.apache.org/download.cgi>

#### Files

Maven is distributed in several formats for your convenience. Simply pick a ready-made binar the [installation instructions](#). Use a source archive if you intend to build Maven yourself.

In order to guard against corrupted downloads/installations, it is highly recommended to [verify](#) bundles against the public [KEYS](#) used by the Apache Maven developers.

	Link	Checksum
Binary tar.gz archive	<a href="#">apache-maven-3.5.2-bin.tar.gz</a>	<a href="#">apache-maven-3.5.2-bin.tar.gz.md5</a>
Binary zip archive	<a href="#">apache-maven-3.5.2-bin.zip</a>	<a href="#">apache-maven-3.5.2-bin.zip.md5</a>
Source tar.gz archive	<a href="#">apache-maven-3.5.2-src.tar.gz</a>	<a href="#">apache-maven-3.5.2-src.tar.gz.md5</a>
Source zip archive	<a href="#">apache-maven-3.5.2-src.zip</a>	<a href="#">apache-maven-3.5.2-src.zip.md5</a>

Extract Maven files.

```
$ cd /home/Spark/Downloads  
$ tar xvf apache-maven-3.5.2-bin.tar.gz
```

Use the following commands to move the Maven files to the directory /usr/local/maven

```
$ su -  
Password:  
# cd /home/Spark/Downloads/  
# mv apache-maven-3.5.2 /usr/local/maven  
# exit
```

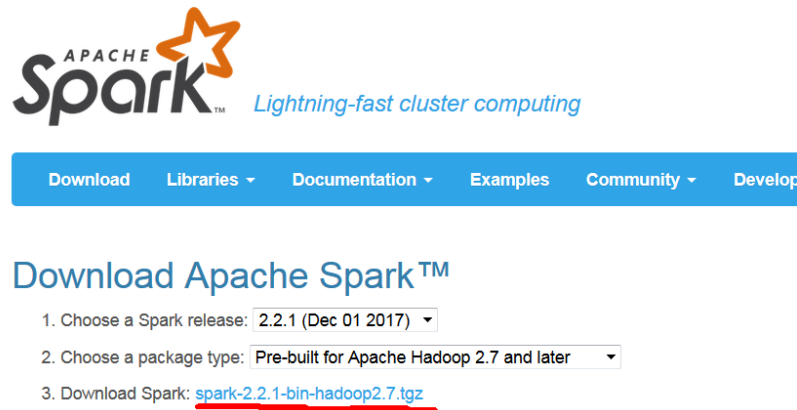
Set the path for Maven.

```
$ sudo gedit /etc/environment
```

Append the following clause at the end of PATH = “/usr/local/sbin:.....” in the file, and save it.  
*:/usr/local/maven/bin*

## 4. Install Spark

Download Spark from <https://spark.apache.org/downloads.html>



Extract the file

```
$ cd /home/Spark/Downloads
$ tar xvf spark-2.2.1-bin-hadoop2.7.tgz
```

Move the Spark files to the directory /usr/local/spark

```
$ su -
Password:
# cd /home/Spark/Downloads/
# mv spark-2.2.1-bin-hadoop2.7 /usr/local/spark
# exit
```

Set the path for Spark.

```
$ sudo gedit /etc/environment
```

Append the following clause at the end of PATH = "/usr/local/sbin:.....", then save it.

*:/usr/local/spark/bin*

**Now, restart the system to make those changes work!**

## 5. Verify the Software Installations

```
$ java -version
```

If Java is installed successfully then you will find the following output.

```
spark@spark-VirtualBox: ~  
spark@spark-VirtualBox:~$ java -version  
java version "1.8.0_151"  
Java(TM) SE Runtime Environment (build 1.8.0_151-b12)  
Java HotSpot(TM) 64-Bit Server VM (build 25.151-b12, mixed mode)  
spark@spark-VirtualBox:~$
```

```
$ scala -version
```

If Scala is installed successfully then you will find the following output.

```
spark@spark-VirtualBox: ~  
spark@spark-VirtualBox:~$ scala -version  
Scala code runner version 2.12.4 -- Copyright 2002-2017, LAMP/EPFL and Lightbend  
, Inc.  
spark@spark-VirtualBox:~$
```

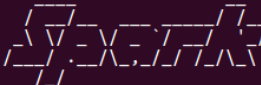
```
$ mvn -version
```

If Maven is installed successfully then you will find the following output.

```
spark@spark-VirtualBox:~$ mvn -version  
Apache Maven 3.5.2 (138ed61fd100ec658bfa2d307c43b76940a5d7d; 2017-10-18T15:58:1  
3+08:00)  
Maven home: /usr/local/maven  
Java version: 1.8.0_151, vendor: Oracle Corporation  
Java home: /usr/lib/jvm/java-8-oracle/jre  
Default locale: en_US, platform encoding: UTF-8  
OS name: "linux", version: "4.10.0-42-generic", arch: "amd64", family: "unix"
```

```
$ spark-shell
```

If spark is installed successfully then you will find the following output.

```
spark@spark-VirtualBox:~$ spark-shell  
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties  
Setting default log level to "WARN".  
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).  
18/01/08 18:59:33 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable  
18/01/08 18:59:33 WARN Utils: Your hostname, spark-VirtualBox resolves to a loop back address: 127.0.0.1; using 10.0.2.15 instead (on interface enp0s3)  
18/01/08 18:59:33 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another address  
Spark context Web UI available at http://10.0.2.15:4040  
Spark context available as 'sc' (master = local[*], app id = local-1515409175769).  
Spark session available as 'spark'.  
Welcome to  
 version 2.2.1  
Using Scala version 2.11.8 (Java HotSpot(TM) 64-Bit Server VM, Java 1.8.0_151)  
Type in expressions to have them evaluated.  
Type :help for more information.  
scala>
```

## Appendix C. My First Spark Program (with Java)

1. Download the example files from IVLE (in.txt, pom.xml, SparkWordCount.java).
2. Create a new folder named "spark-application" with the following files.

```
.
./pom.xml
./in.txt
./src
./src/main
./src/main/java
./src/main/java/com
./src/main/java/com/mycompany
./src/main/java/com/mycompany/app
./src/main/java/com/mycompany/app/SparkWordCount.java
```

3. Execute the following command under "." directory to generate the .jar file.

```
$ mvn package
```

```
spark@spark-VirtualBox:~/Downloads/spark-java/spark-application$ mvn package
[INFO] Scanning for projects...
[INFO]
[INFO] -----
[INFO] Building wordcount 1.0
[INFO] -----
[INFO]
[INFO] --- maven-resources-plugin:2.6:resources (default-resources) @ wordcount
```

...

```
[INFO] Building jar: /home/spark/Downloads/spark-java/spark-application/target/w
ordcount-1.0.jar
[INFO] -----
[INFO] BUILD SUCCESS
[INFO] -----
[INFO] Total time: 4.203 s
[INFO] Finished at: 2018-01-09T15:34:16+08:00
[INFO] Final Memory: 33M/79M
[INFO] -----
```

There will be a new folder named *target* generated in "." directory.



4. Submit to Spark and execute the program using the following command.

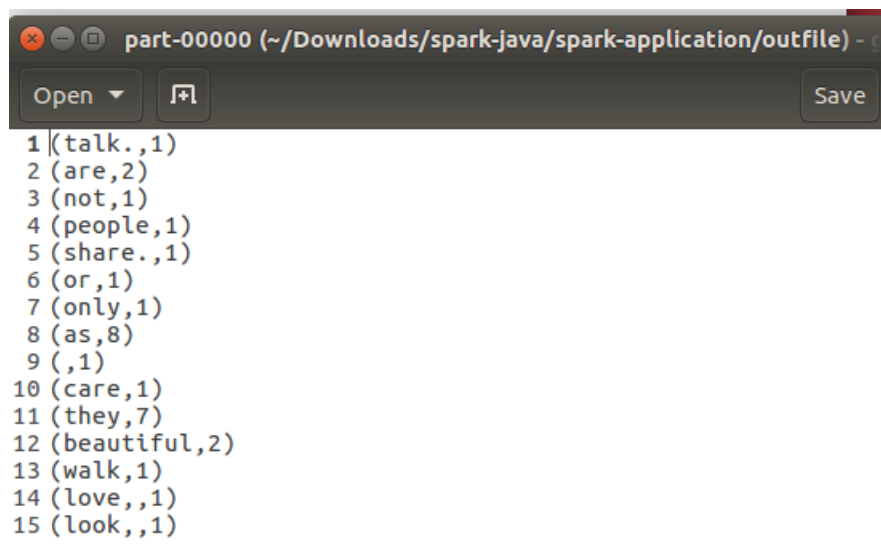
```
$ spark-submit --class SparkWordCount --master local target/wordcount-1.0.jar
```

```
spark@spark-VirtualBox:~/Downloads/spark-java/spark-application$ spark-submit --
class SparkWordCount --master local target/wordcount-1.0.jar
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties
18/01/09 15:37:20 INFO SparkContext: Running Spark version 2.2.1
18/01/09 15:37:20 WARN NativeCodeLoader: Unable to load native-hadoop library fo
r your platform... using builtin-java classes where applicable
18/01/09 15:37:21 WARN Utils: Your hostname, spark-VirtualBox resolves to a loop
back address: 127.0.1.1; using 10.0.2.15 instead (on interface enp0s3)
18/01/09 15:37:21 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another
address
18/01/09 15:37:21 INFO SparkContext: Submitted application: wordcount
18/01/09 15:37:21 INFO SecurityManager: Changing view acls to: spark
18/01/09 15:37:21 INFO SecurityManager: Changing modify acls to: spark
18/01/09 15:37:21 INFO SecurityManager: Changing view acls groups to:
18/01/09 15:37:21 INFO SecurityManager: Changing modify acls groups to:
18/01/09 15:37:21 INFO SecurityManager: SecurityManager: authentication disabled
; ui acls disabled; users with view permissions: Set(spark); groups with view p
ermissions: Set(); users with modify permissions: Set(spark); groups with modif
y permissions: Set()
```

...

```
18/01/09 15:38:54 INFO MemoryStore: MemoryStore cleared
18/01/09 15:38:54 INFO BlockManager: BlockManager stopped
18/01/09 15:38:54 INFO BlockManagerMaster: BlockManagerMaster stopped
18/01/09 15:38:54 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint:
OutputCommitCoordinator stopped!
18/01/09 15:38:54 INFO SparkContext: Successfully stopped SparkContext
18/01/09 15:38:54 INFO ShutdownHookManager: Shutdown hook called
18/01/09 15:38:54 INFO ShutdownHookManager: Deleting directory /tmp/spark-fe3e56
8d-1d5f-4d7a-af40-ff2798cf1b98
spark@spark-VirtualBox:~/Downloads/spark-java/spark-application$
```

You can see a folder named *outfile* generated under "." directory. The result is in the inside file named *part-00000*.



The screenshot shows a file explorer window titled "part-00000 (~/.Downloads/spark-java/spark-application/outfile)". The window contains a list of 15 lines of text, each representing a word and its frequency in parentheses. The text is as follows:

```
1 (talk.,1)
2 (are,2)
3 (not,1)
4 (people,1)
5 (share.,1)
6 (or,1)
7 (only,1)
8 (as,8)
9 (,1)
10 (care,1)
11 (they,7)
12 (beautiful,2)
13 (walk,1)
14 (love,,1)
15 (look,,1)
```

## Appendix D. My First Spark Program (with Python)

1. Download the example files from IVLE (in.txt, wordcount.py).
2. Create a new folder named "spark-application" with the files you downloaded. (in.txt, wordcount.py).
3. To execute the Spark program, using the following command under the folder ".../spark-application/".

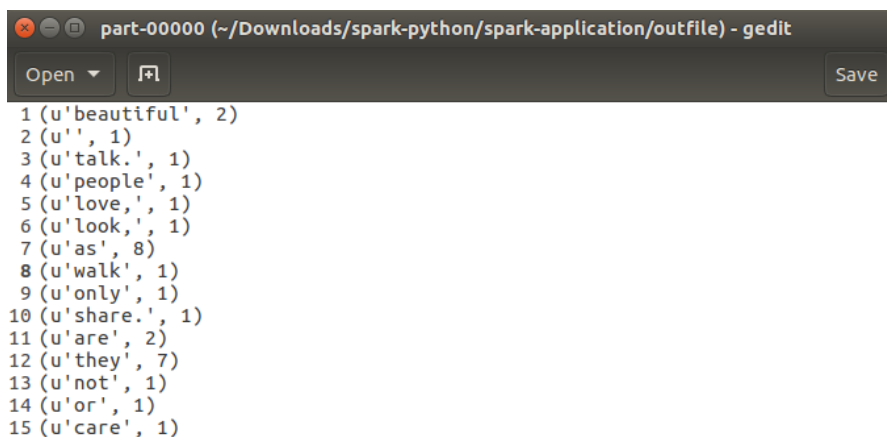
```
$ spark-submit wordcount.py in.txt outfile
```

```
spark@spark-VirtualBox:~/Downloads/spark-python/spark-application$ spark-submit
wordcount.py in.txt outfile
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties
19/01/14 12:13:02 WARN Utils: Your hostname, spark-VirtualBox resolves to a loop
back address: 127.0.1.1; using 10.0.2.15 instead (on interface enp0s3)
19/01/14 12:13:02 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another
address
19/01/14 12:13:04 INFO SparkContext: Running Spark version 2.2.1
19/01/14 12:13:05 WARN NativeCodeLoader: Unable to load native-hadoop library fo
r your platform... using builtin-java classes where applicable
19/01/14 12:13:05 INFO SparkContext: Submitted application: wordcount.py
19/01/14 12:13:05 INFO SecurityManager: Changing view acls to: spark
19/01/14 12:13:05 INFO SecurityManager: Changing modify acls to: spark
```

...

```
19/01/14 12:13:12 INFO SparkContext: Successfully stopped SparkContext
19/01/14 12:13:12 INFO ShutdownHookManager: Shutdown hook called
19/01/14 12:13:12 INFO ShutdownHookManager: Deleting directory /tmp/spark-cd9294
59-3369-4305-af19-9427fde05ea3
19/01/14 12:13:12 INFO ShutdownHookManager: Deleting directory /tmp/spark-cd9294
59-3369-4305-af19-9427fde05ea3/pyspark-4dfe6e40-097f-484d-8900-72346cfb353f
spark@spark-VirtualBox:~/Downloads/spark-python/spark-application$
```

You can see a folder named *outfile* generated under "." directory. The result is in the inside file named *part-00000*.



```
part-00000 (~/Downloads/spark-python/spark-application/outfile) - gedit
Open Save
1 (u'beautiful', 2)
2 (u'', 1)
3 (u'talk.', 1)
4 (u'people', 1)
5 (u'love,', 1)
6 (u'look,', 1)
7 (u'as', 8)
8 (u'walk', 1)
9 (u'only', 1)
10 (u'share.', 1)
11 (u'are', 2)
12 (u'they', 7)
13 (u'not', 1)
14 (u'or', 1)
15 (u'care', 1)
```