

# Contents

<b>Exercise 1: Setting up the Environment</b> .....	2
<b>Exercise 2: Practicing HDFS Commands</b> .....	4
<b>Exercise 3: Running MapReduce Job</b> .....	8
<b>Exercise 4: Launching Spark</b> .....	23
<b>Exploring Data using RDD operations:</b> .....	25
Transformations.....	25
Action.....	30
PairRDD.....	34
<b>Developing with Spark</b> .....	40
REPL.....	40
Zeppelin.....	42
<b>Exercise 5: Building Spark Application using Eclipse IDE</b> .....	43
<b>Additional Exercise</b> .....	52
Analyze Movie Lens Dataset using RDD.....	52
<b>Exercise 6: Dataframe and Spark SQL</b> .....	58
Infer schema by Reflection (case class) .....	58
Programmatically .....	63
UDF in Spark Dataframe.....	65
<b>Exercise 7: Spark Streaming</b> .....	68
<b>Exercise 8: SparkML</b> .....	70
Predicting power grid demand using Spark ML .....	70
<b>Exercise9: Machine Learning</b> .....	85
Linear Regression .....	86
Logistic Regression .....	98
Decision Tree.....	99
Random Forest .....	102
Classification.....	104
Clustering .....	113
PCA (Principle Component Analysis) .....	125
<b>Exercise10: Graphx</b> .....	129
Analyzing Flight Data.....	129

## Exercise 1: Setting up the Environment

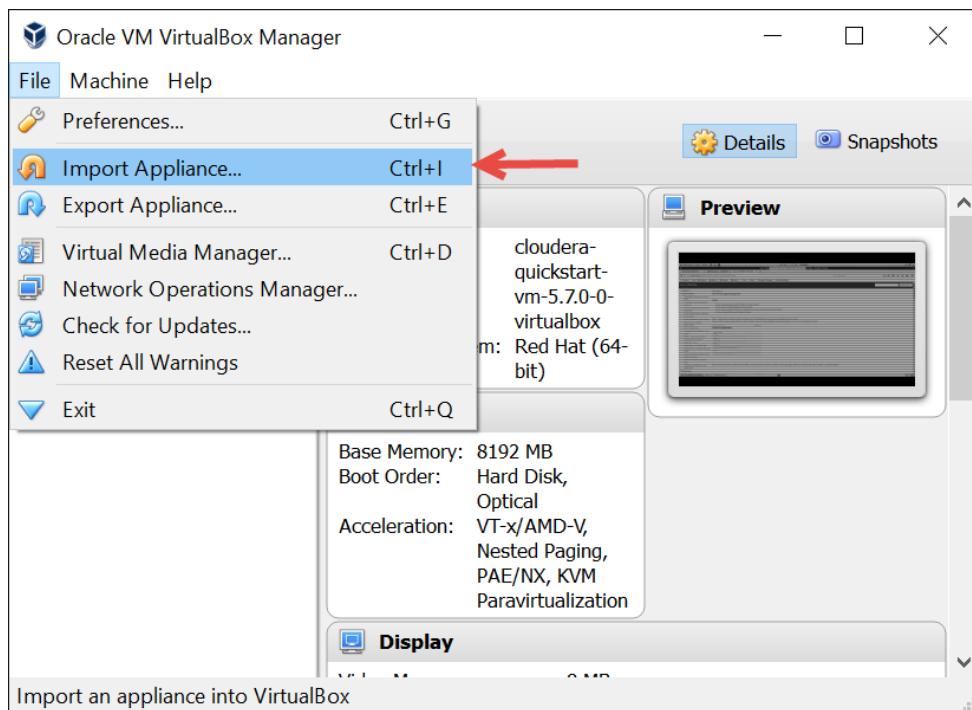
Step1: Installing Oracle VM virtual box from Training Bundle.

- Under ‘Training Bundle’ folder, navigate to “Software” folder.
- Find the executable file named ‘VirtualBox-5.0.16-105871-Win’ and complete the installation.

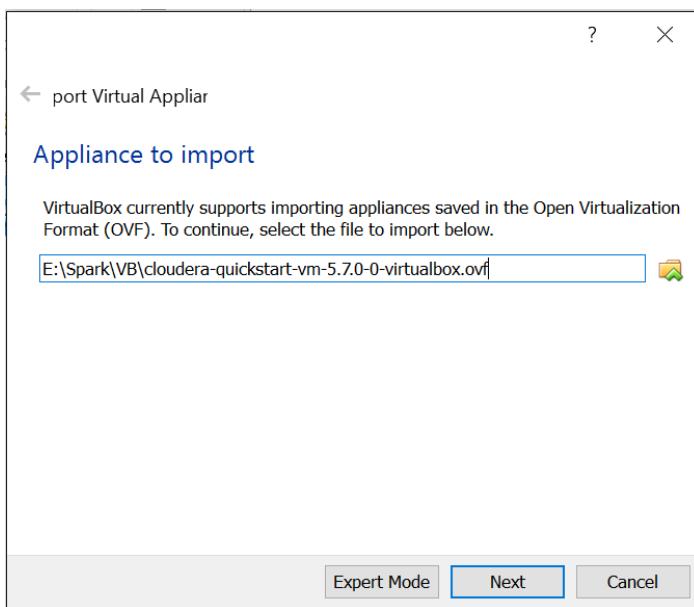
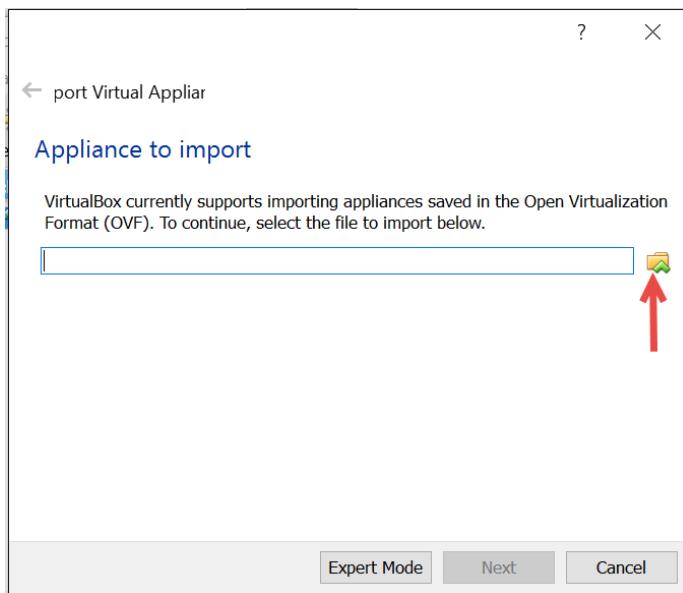
Step2: Extracting the Image

- From the same training bundle, locate the folder named VM image> ‘cloudera-quickstart-vm-5.7.0-0-virtualbox’
- Load the image to get started with the machine.

Step3: Importing the file



- Browse to the location under “Training Bundle” and select the OVF file.



- This will prepare your machine.
- After import, change the number of processors to two before starting the machine.
  
- Verify all the running services with jps command

```
[cloudera@quickstart ~]$ sudo jps
6403 RESTServer
5709 SecondaryNameNode
5980 NodeManager
6878 RunJar
8543
10329 ZeppelinServer
9721 org.eclipse.equinox.launcher_1.3.0.v20140415-2008.jar
5841 Bootstrap
5366 DataNode
5454 JournalNode
8584
7516 HistoryServer
8518 Bootstrap
5573 NameNode
7484 Bootstrap
5303 QuorumPeerMain
6659 RunJar
5897 JobHistoryServer
18565 Jps
6184 ResourceManager
6521 ThriftServer
8191 Bootstrap
[cloudera@quickstart ~]$
```

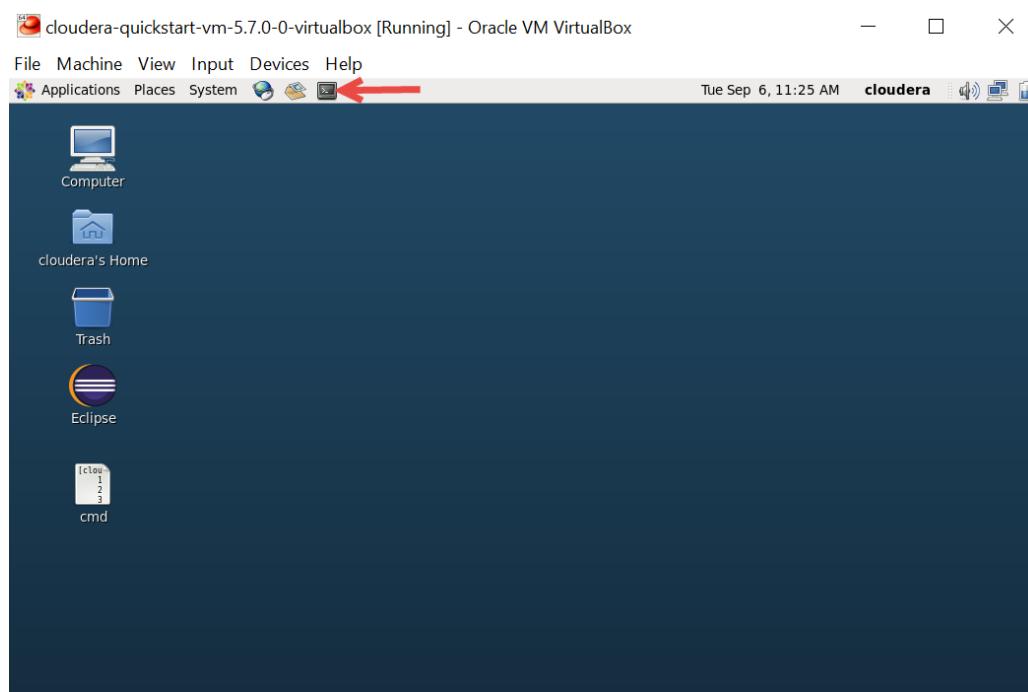
## Exercise 2: Practicing HDFS Commands

### Task 1: Using Hadoop Command Line Interface

Navigate to Application > System Tools > Terminal

OR

Use shortcut to open the terminal.



- Operation on Files and directories

To check the content of root directory in HDFS

```
hdfs dfs -ls /
```

```
[cloudera@quickstart ~]$ hdfs dfs -ls /
Found 5 items
drwxrwxrwx  - hdfs  supergroup      0 2016-04-05 23:56 /benchmarks
drwxr-xr-x  - hbase supergroup      0 2016-09-04 10:14 /hbase
drwxrwxrwt  - hdfs  supergroup      0 2016-08-03 01:56 /tmp
drwxr-xr-x  - hdfs  supergroup      0 2016-08-02 08:08 /user
drwxr-xr-x  - hdfs  supergroup      0 2016-04-05 23:58 /var
[cloudera@quickstart ~]$ █
```

View the content of /user directory

```
hdfs dfs -ls /user
```

```
[cloudera@quickstart ~]$ hdfs dfs -ls /user
Found 9 items
drwxr-xr-x  - cloudera cloudera      0 2016-08-03 01:54 /user/cloudera
drwxr-xr-x  - hdfs    supergroup      0 2016-08-02 08:08 /user/hdfs
drwxr-xr-x  - mapred   hadoop        0 2016-04-05 23:56 /user/history
drwxrwxrwx  - hive     supergroup      0 2016-04-05 23:58 /user/hive
drwxrwxrwx  - hue      supergroup      0 2016-04-05 23:57 /user/hue
drwxrwxrwx  - jenkins  supergroup      0 2016-04-05 23:56 /user/jenkins
drwxrwxrwx  - oozie    supergroup      0 2016-04-05 23:57 /user/oozie
drwxrwxrwx  - root     supergroup      0 2016-04-05 23:57 /user/root
drwxr-xr-x  - hdfs    supergroup      0 2016-04-05 23:58 /user/spark
```



For an empty directory the prompt does not show any error while querying whereas if the directory doesn't exist it will throw an error.

Example: Query a non-existing directory say /music

```
hdfs dfs -ls /music
```

```
[cloudera@quickstart ~]$ hdfs dfs -ls /music
ls: `/music': No such file or directory
[cloudera@quickstart ~]$ █
```



The directory structure of Hadoop filesystem is different from local filesystem i.e. Linux filesystem. To know the difference, list the files on your local filesystem.

Create a new directory named “hadoop” in HDFS

```
hdfs dfs -mkdir hadoop
```

Create a sample.txt file on local and upload the file into newly created directory

```
hdfs dfs -put sample.txt hadoop/sample.txt
```

View the content of new file

```
hdfs dfs -cat hadoop/sample.txt
```



In HDFS, any non-absolute path is considered relative to your home directory i.e. /user/cloudera. Unlike Linux filesystem concept of “current” or “present working directory”.

Download a file from HDFS to your local filesystem, specify HDFS path and local path to achieve the same.

```
hdfs dfs -get /user/cloudera/sample.txt /home/cloudera
```

Remove the directory along with its content (perform this step after completing Task 2)

```
hdfs dfs -rm -r hadoop
```

List all the shell commands supported by Hadoop

```
hdfs dfs
```

## Task2: Using HUE File browser

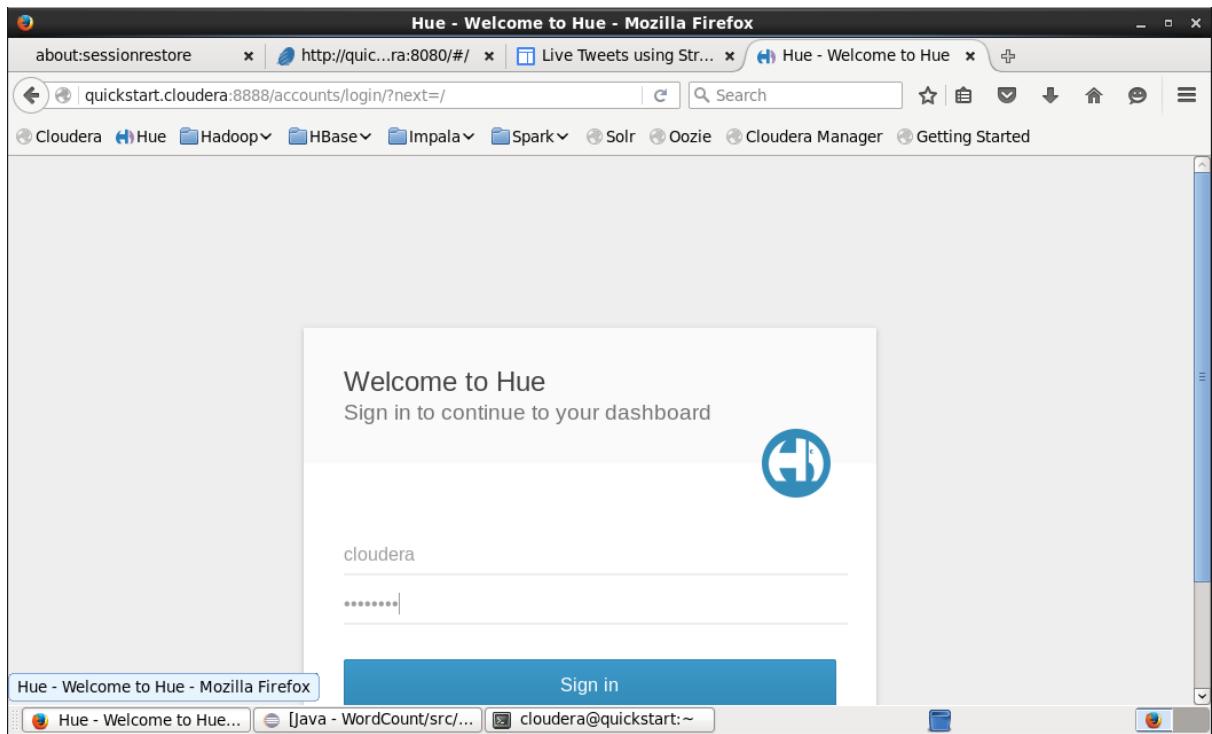
1. Open a web browser on your VM, and navigate to bookmark tab and select **HUE**. This can be alternatively launched by specifying <http://quickstart.cloudera:8888>
2. Enter username as ‘cloudera’ and password ‘cloudera’.
3. To access HDFS, click on **Manage HDFS** in the HUE menu bar
4. By default, the content of home directory i.e. /user/cloudera will be visible.
5. Point out to the directory named “hadoop” created above and view the file ‘sample.txt’
6. Click on **Upload button** and select the appropriate format of the file to be uploaded i.e. plain file or zipped file



The zipped file will be automatically unzipped after upload.

7. Select **Files> Select Files**, and browser to location of data file on your local filesystem

8. Choose the file and click the **Open button**.
9. The selected file will be displayed in the directory /user/cloudera/
10. Click on the checkbox next to file's icon and then tab on **Actions** button to know the possible actions that can be performed on the file.
11. Once you have practised above steps, you can remove the files by clicking on **“Move to Trash”** button.

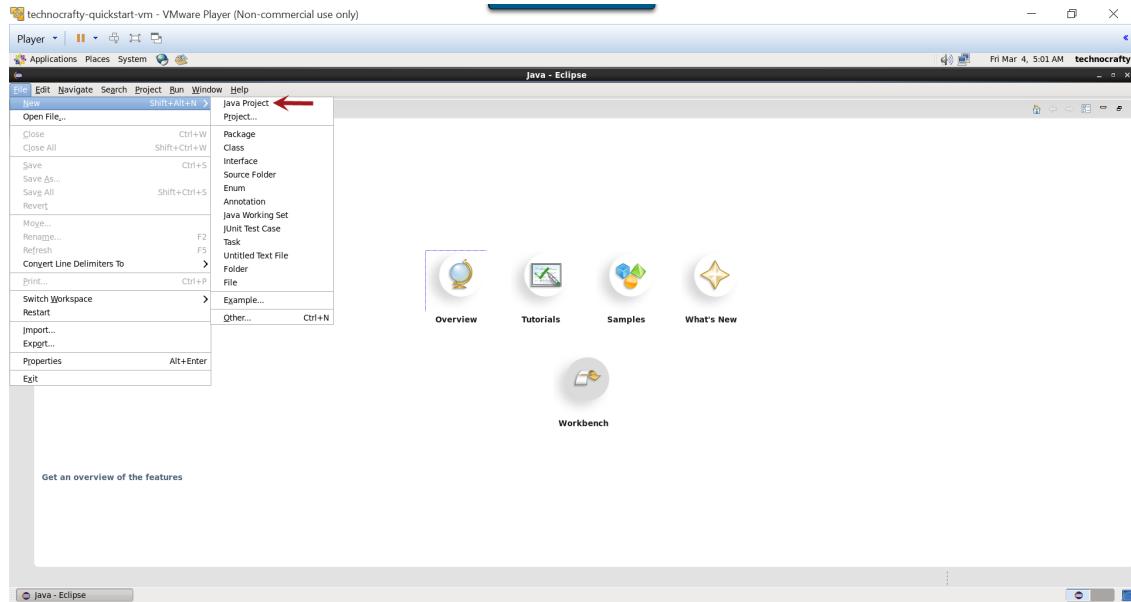


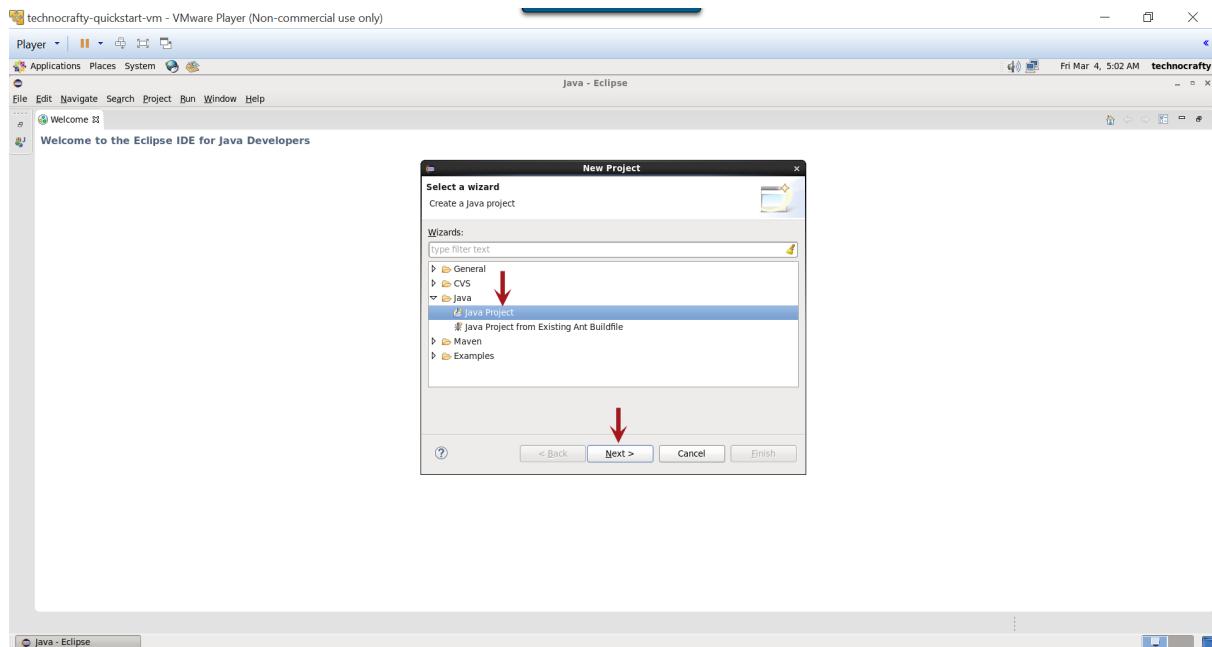
## Exercise 3: Running MapReduce Job

**Prerequisites:** Create an Eclipse WordCount Project

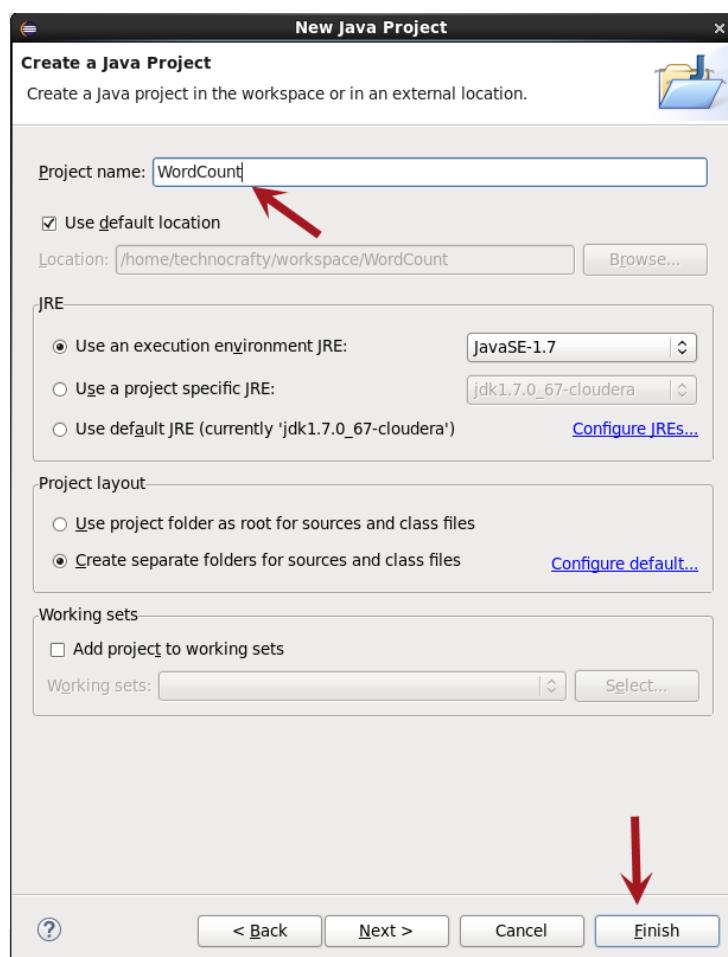
**Step1:** Launch Eclipse, Select the workspace location, and click on “OK”.

**Step2:** Select **File** > **New** > **Java Project** > **Next**

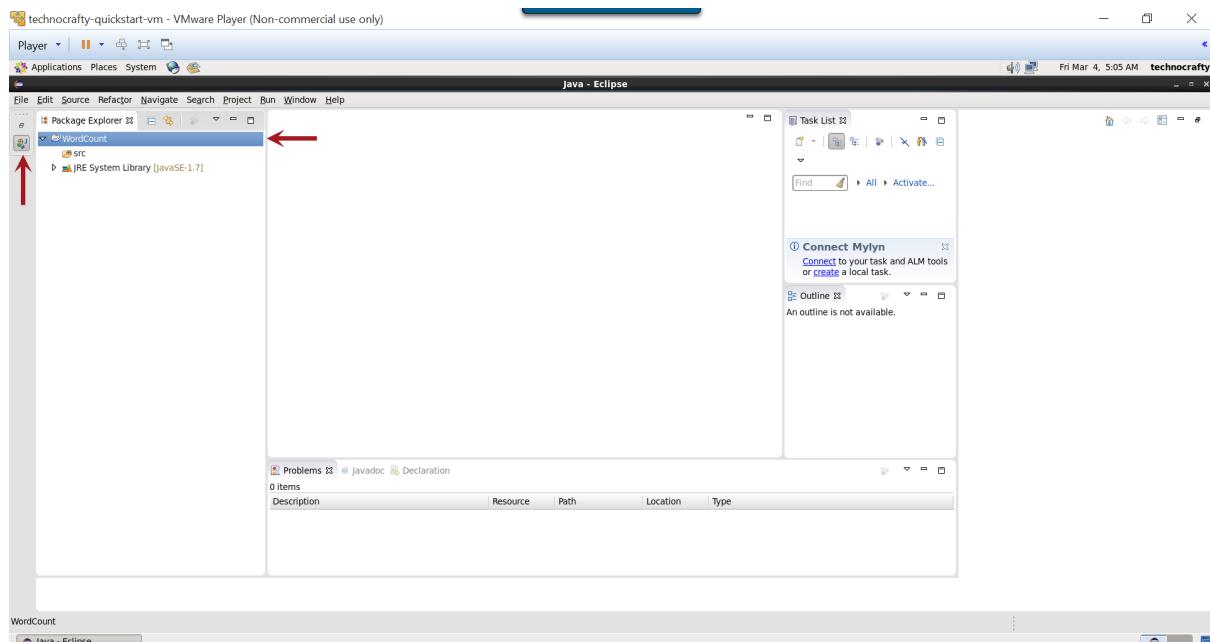




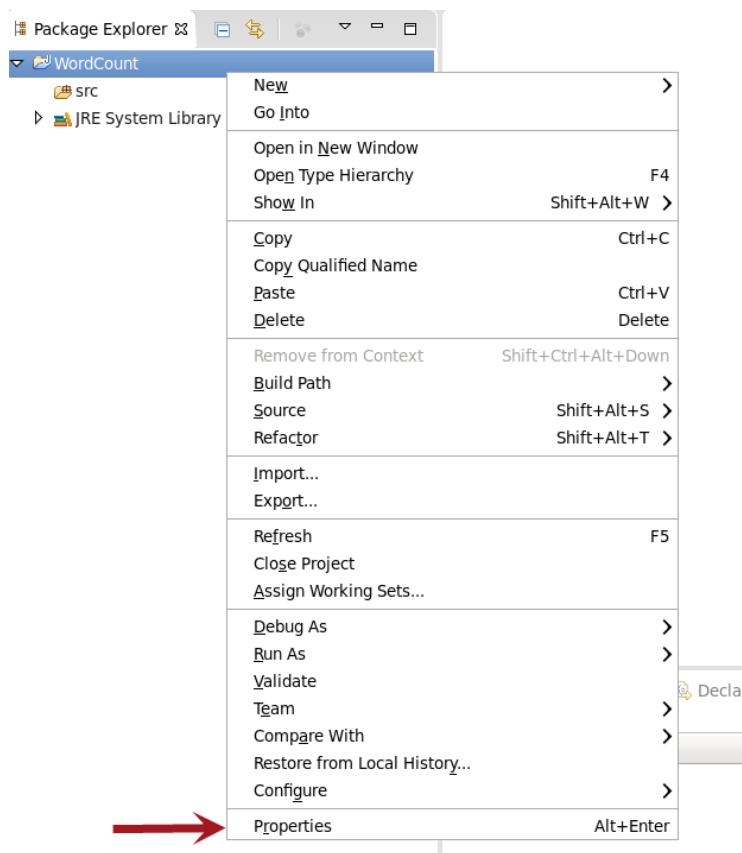
**Step3:** Name the project as "Wordcount" and click "**Finish**"



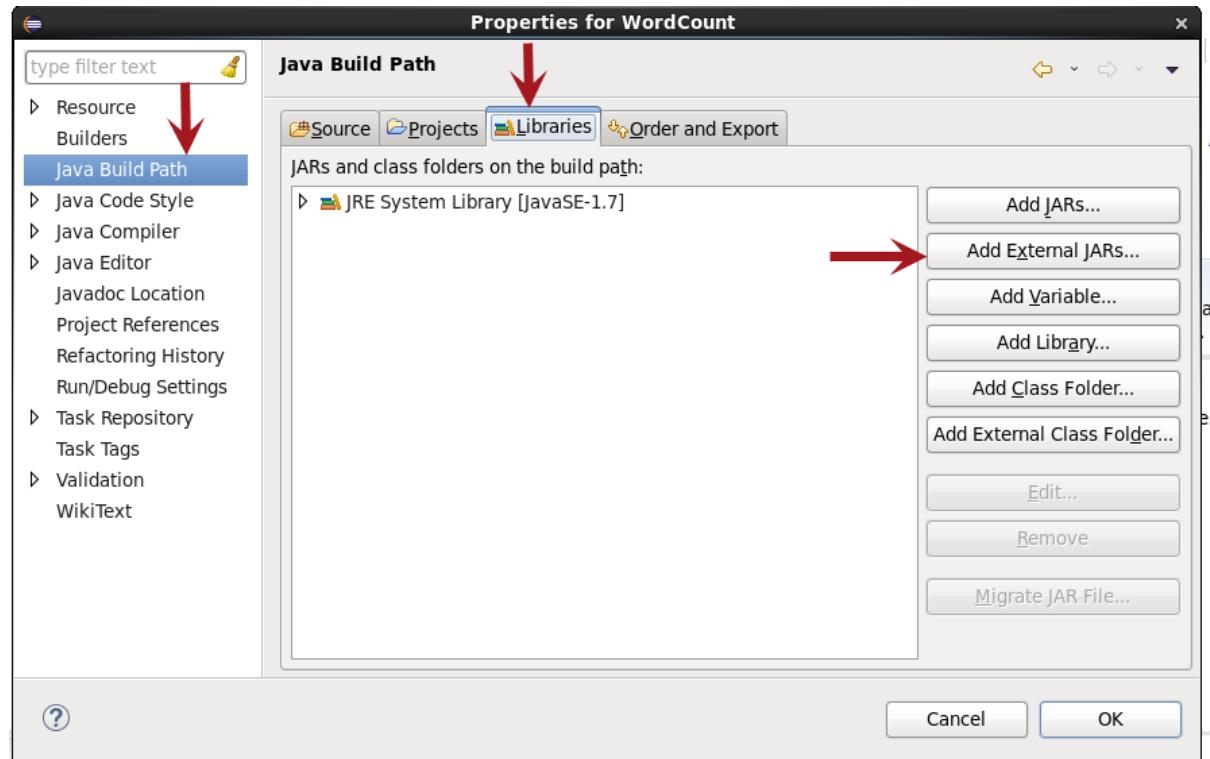
**Step4:** Expand the “Package Explorer” tab and locate your new project i.e. “WordCount”.



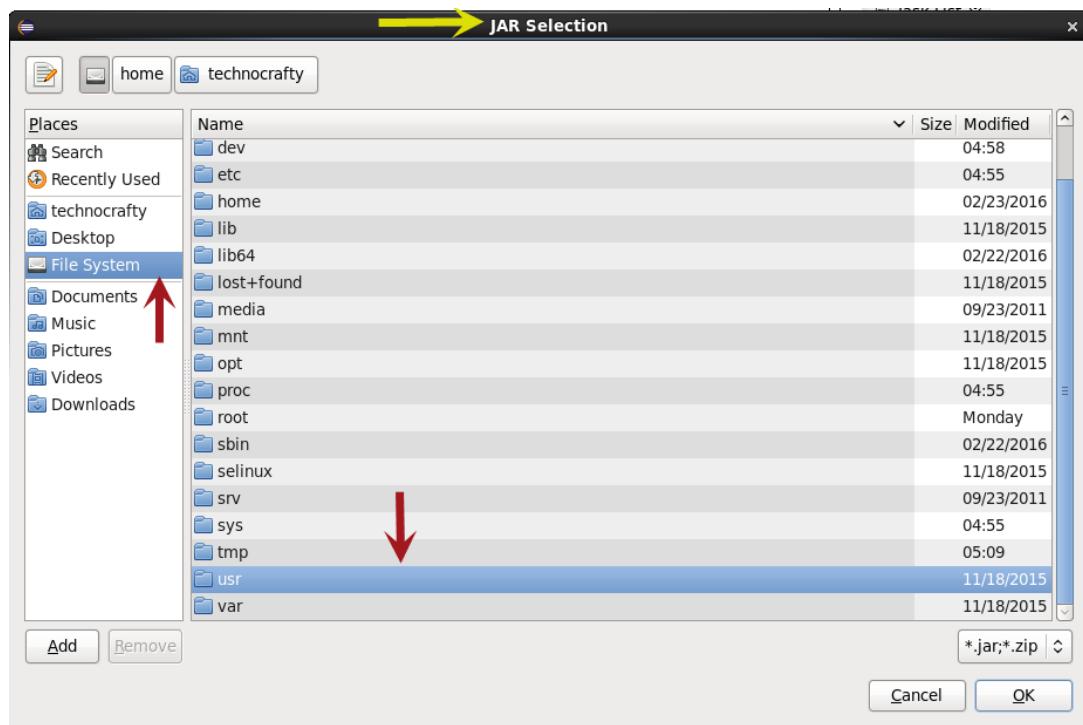
## Step5: Export Hadoop Libraries by Right Click on Wordcount project and selecting Properties

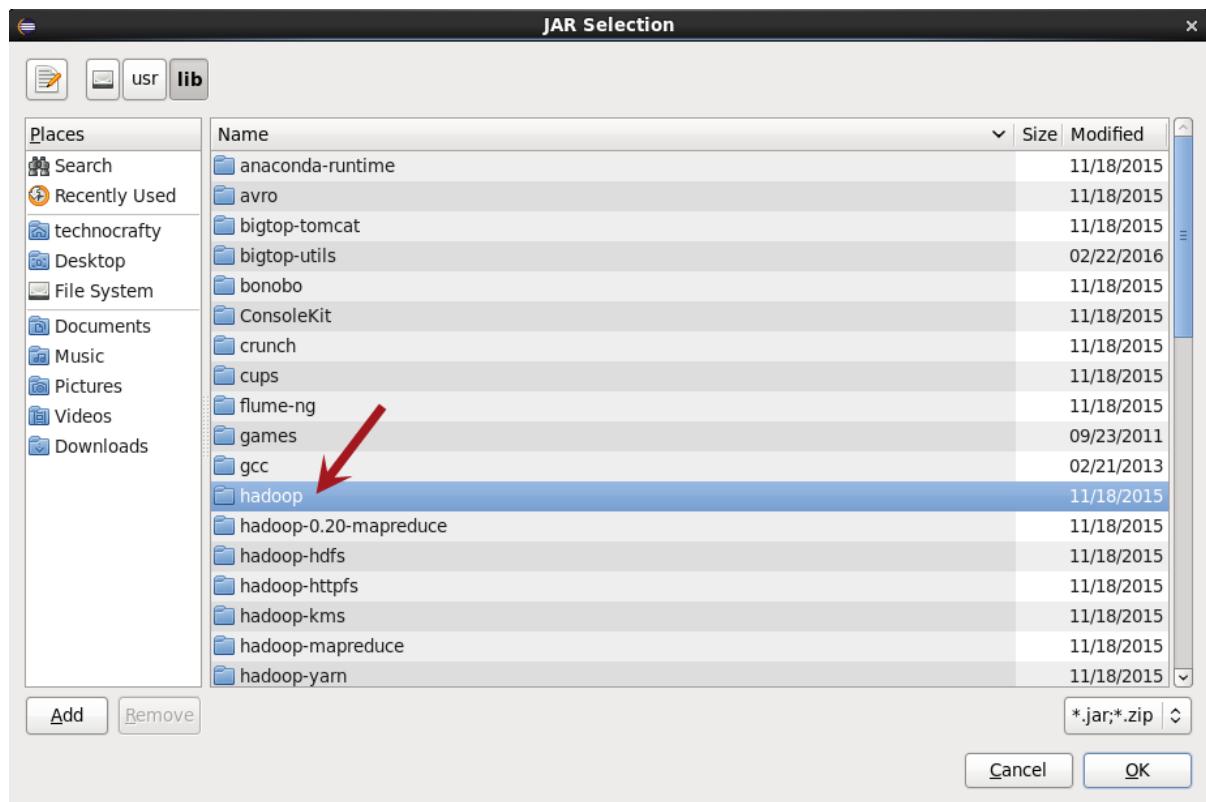
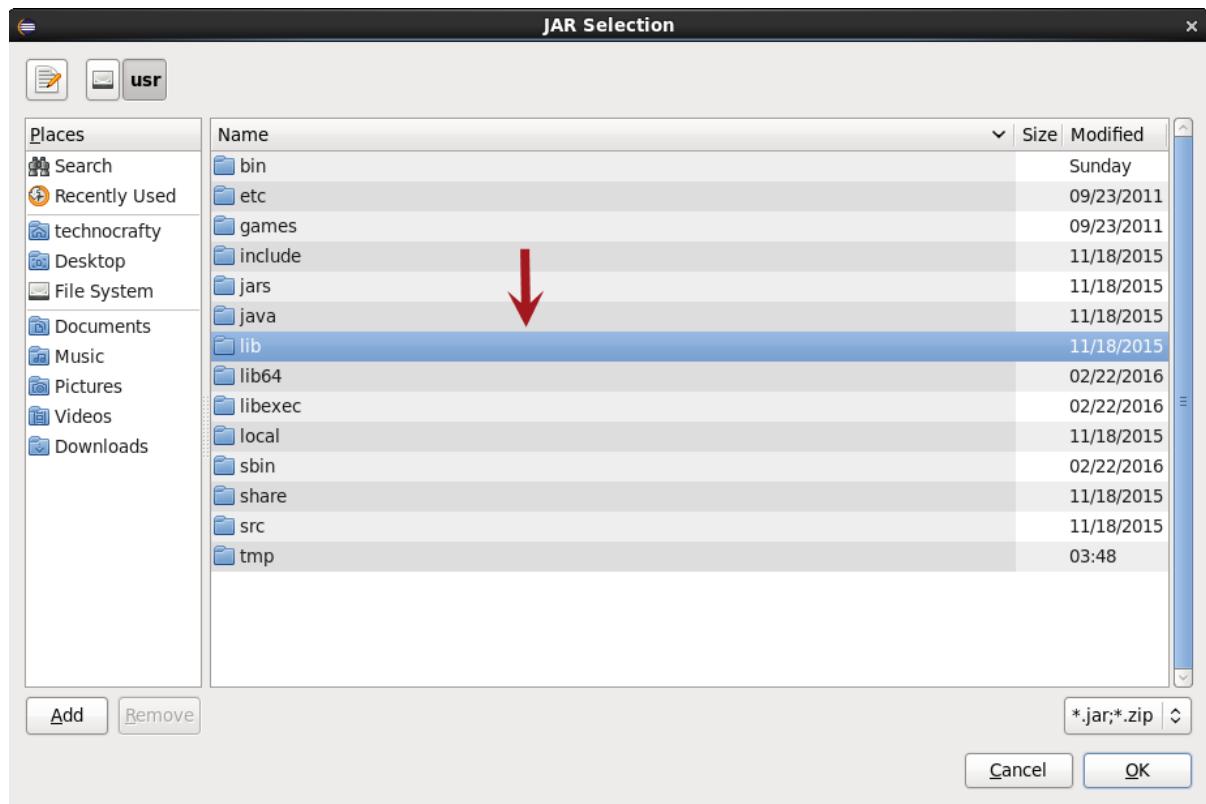


**Step6:** Click on Java Build Path from left panel and select “Libraries” tab from right panel view

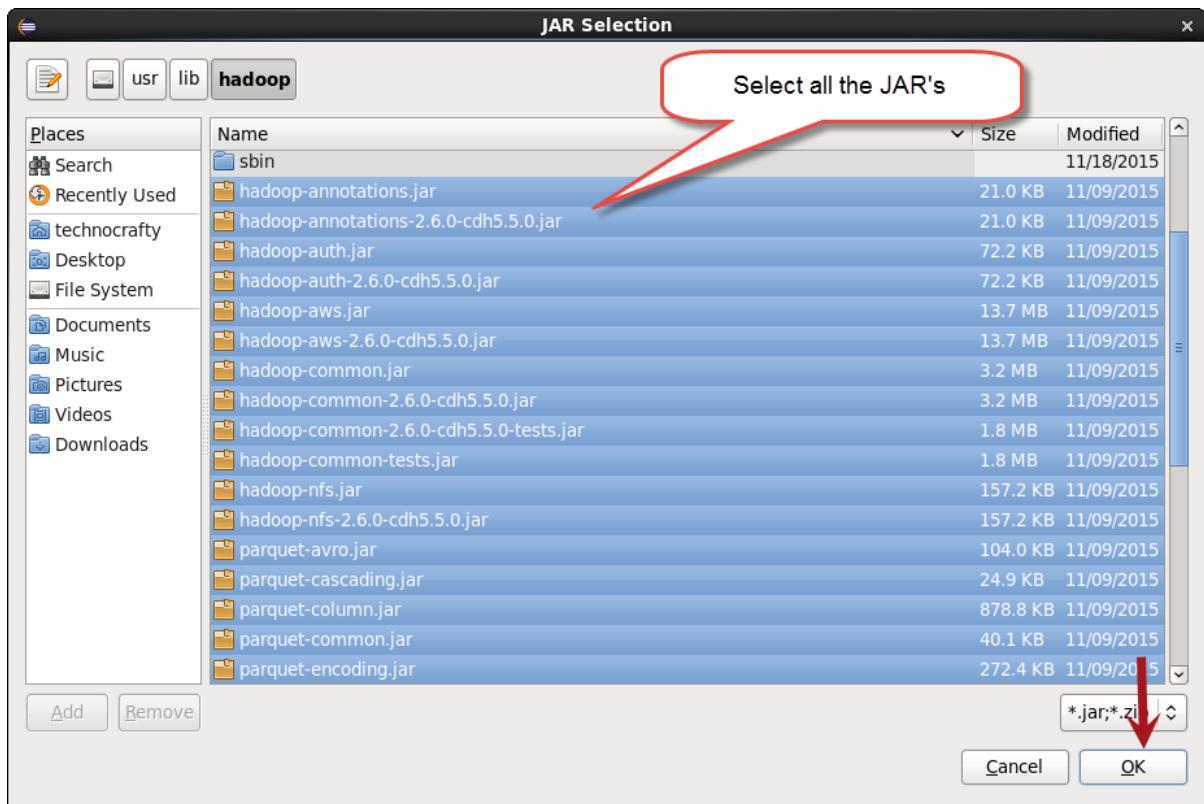


**Step7:** Select “Add External JAR...”, then navigate to File System > usr > lib > Hadoop

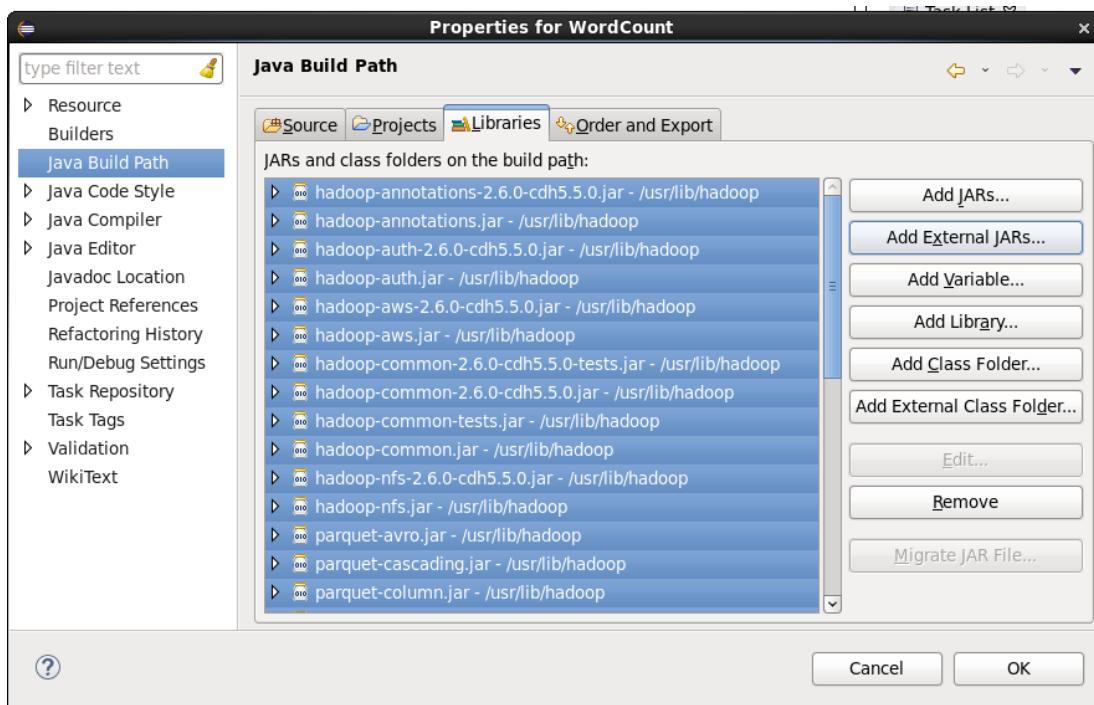




**Step8:** Select all the JAR available under this folder

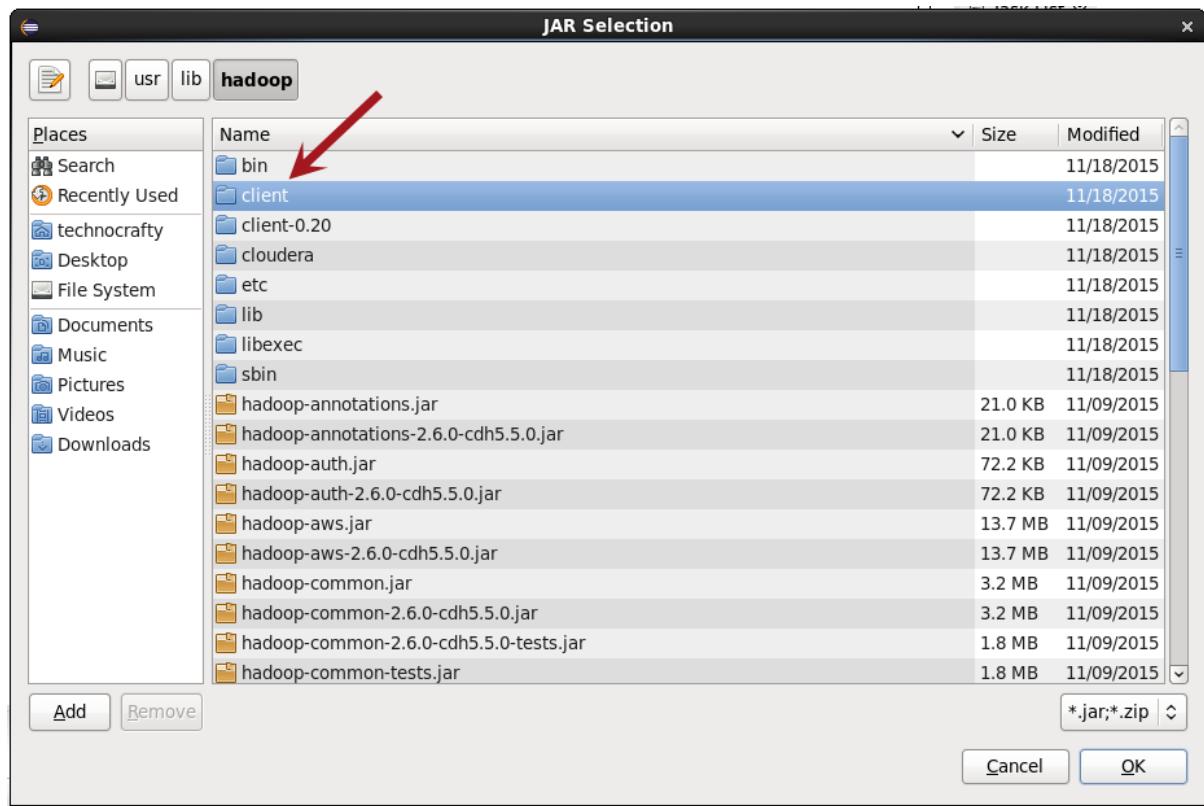


Click on "OK"

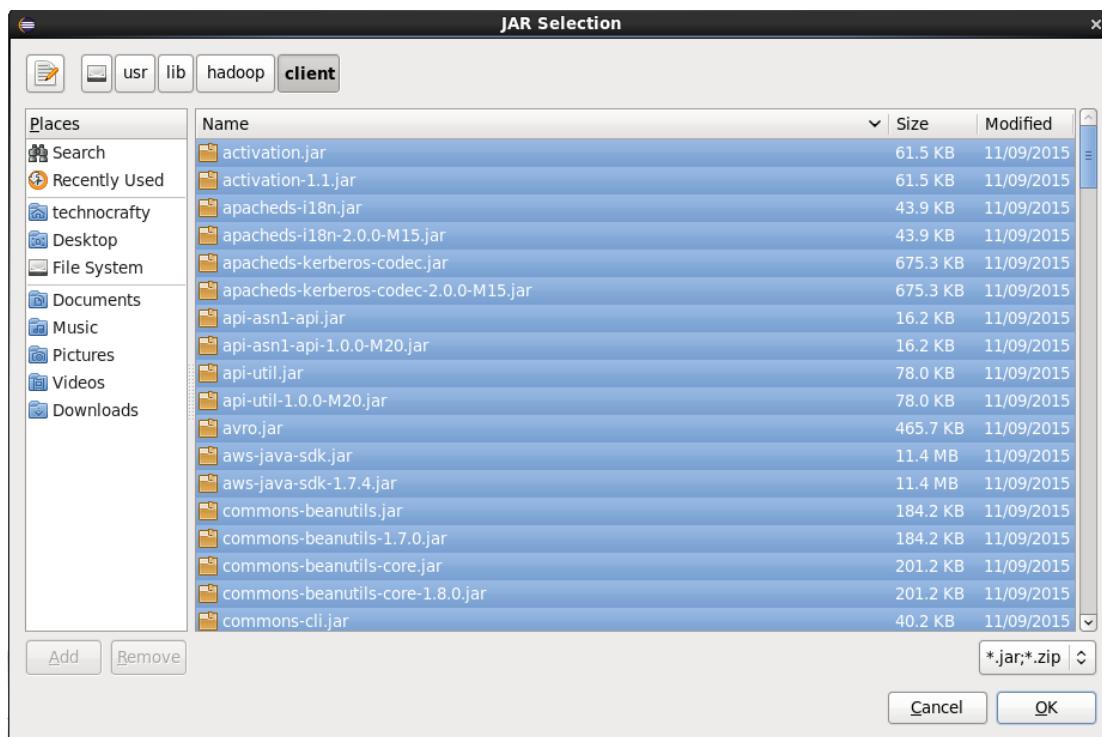


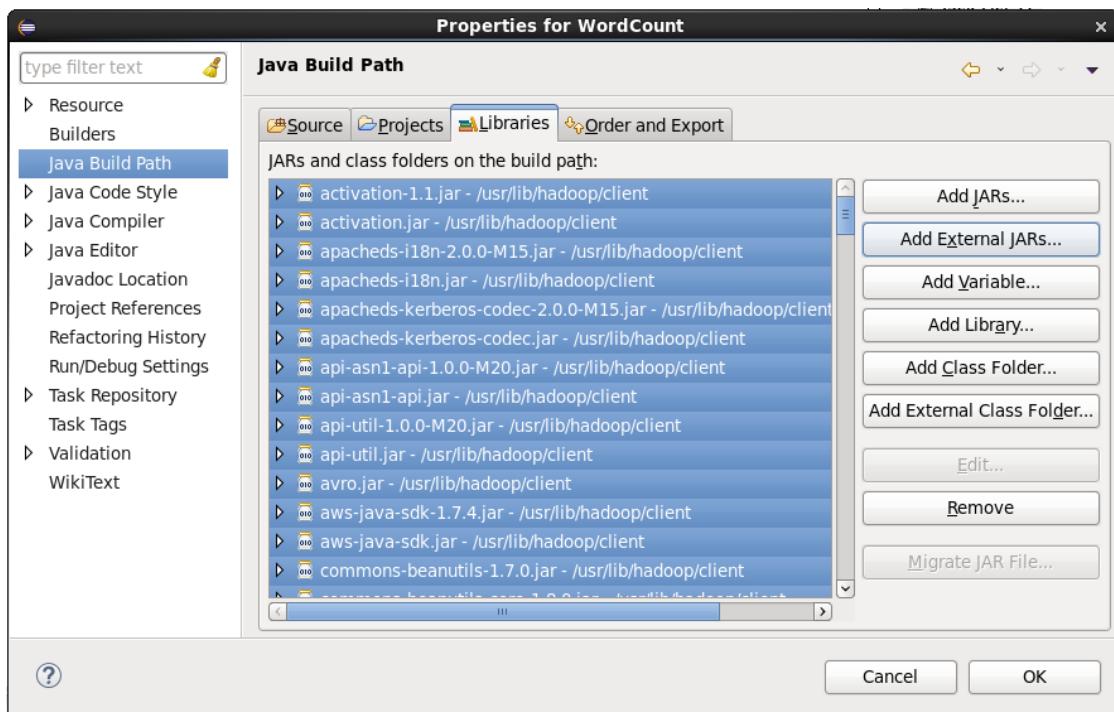
**Step 9:** Perform same steps for /usr/lib/hadoop-mapreduce and /usr/lib/Hadoop-yarn

**Step10:** Now grab all the library JAR files under "Client", by following the same method.

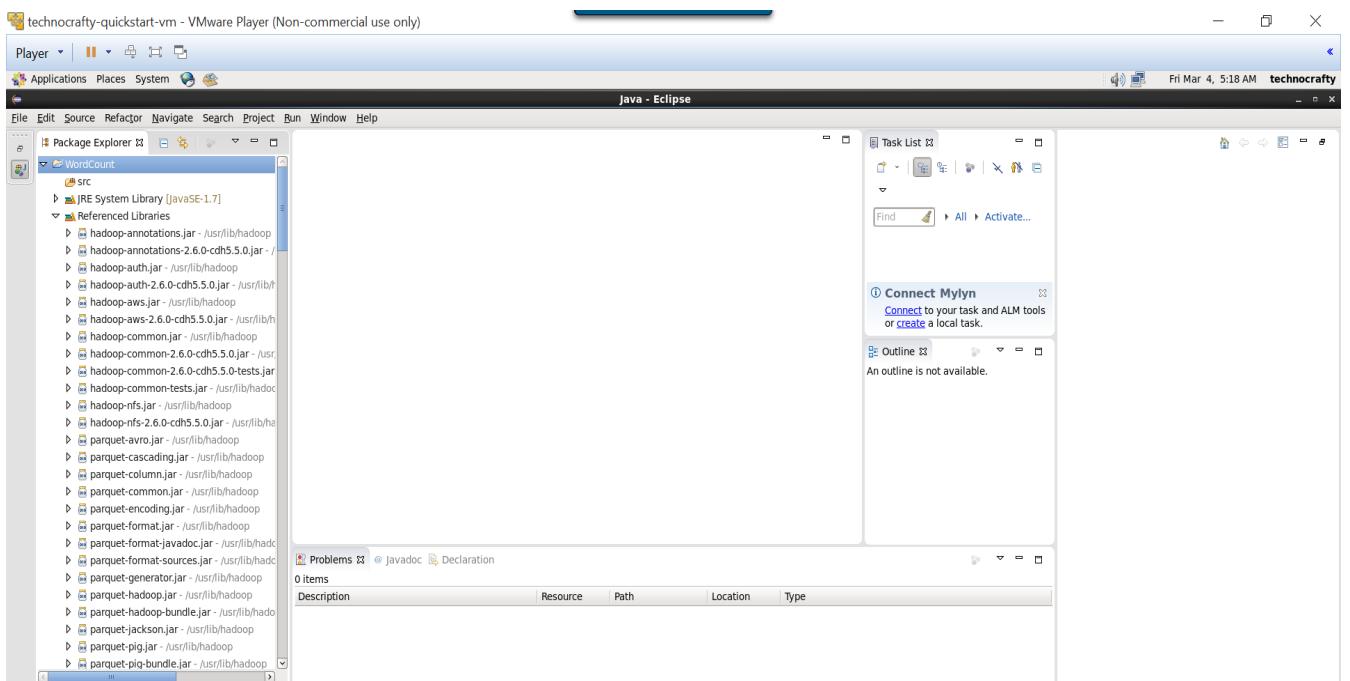


Select the files and click **OK**

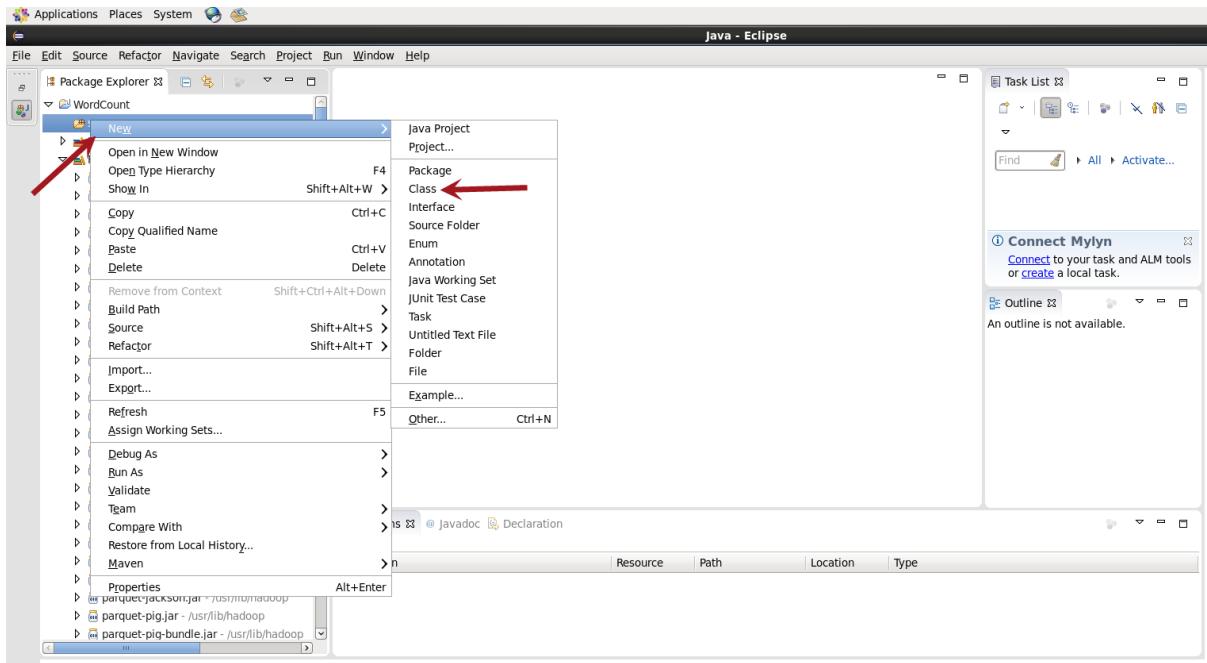




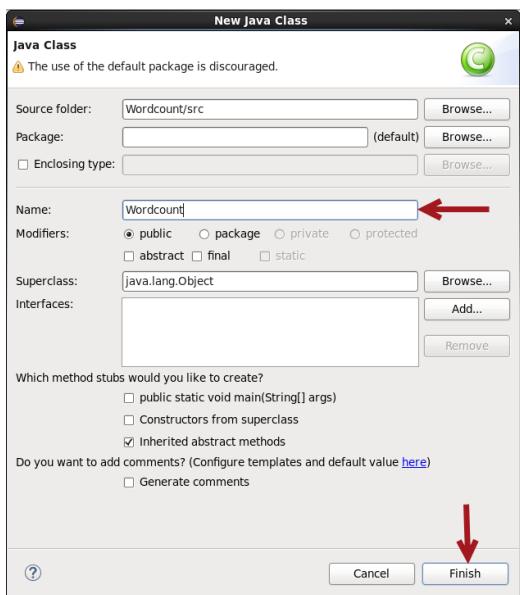
### Step11: Finally, you will find all the JAR files under “Referenced Libraries”



**Step12:** Creating the class files, Right click on src under WordCount and select **New > Class**



**Step13:** Enter the name of Java Class as "Wordcount" and click on **Finish**.



**Step14:** Now it's ready to take your input Mapreduce program

```

1 public class WordCount {
2
3 }

```

## Step15: MapReduce Wordcount Program



## Step 16: Save the program

```

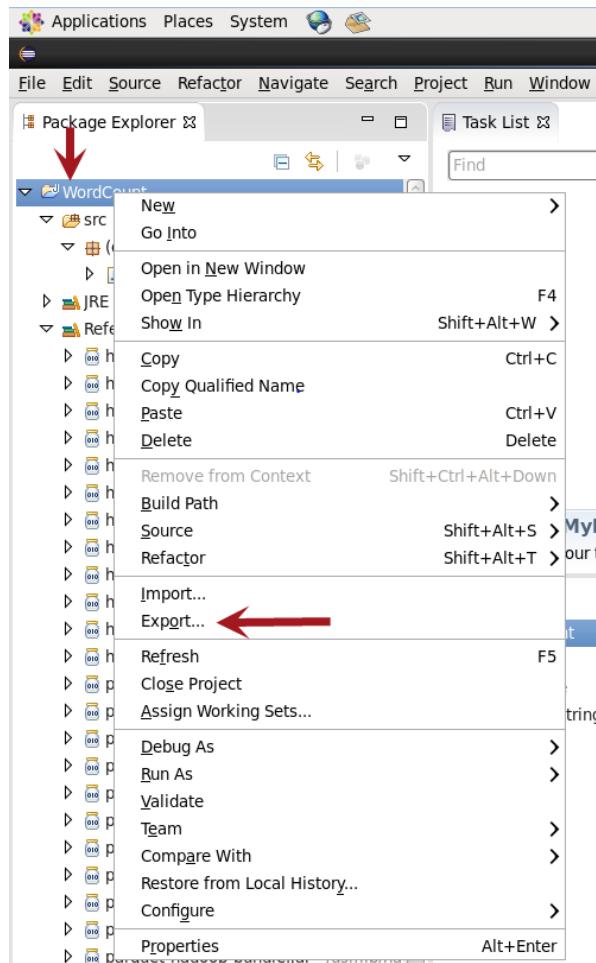
1 import java.io.IOException;
2
3 public class WordCount extends Configured implements Tool {
4
5     public static void main(String args[]) throws Exception {
6         int res = ToolRunner.run(new WordCount(), args);
7         System.exit(res);
8     }
9
10    public int run(String[] args) throws Exception {
11        Path inputPath = new Path(args[0]);
12        Path outputPath = new Path(args[1]);
13
14        Configuration conf = getConf();
15        @SuppressWarnings("deprecation")
16        Job job = new Job(conf, this.getClass().toString());
17
18        FileInputFormat.setInputPaths(job, inputPath);
19        FileOutputFormat.setOutputPath(job, outputPath);
20
21        job.setJobName("WordCount");
22        job.setJarByClass(WordCount.class);
23        job.setInputFormatClass(TextInputFormat.class);
24        job.setOutputFormatClass(TextOutputFormat.class);
25        job.setMapOutputKeyClass(Text.class);
26        job.setMapOutputValueClass(IntWritable.class);
27        job.setOutputKeyClass(Text.class);
28        job.setOutputValueClass(IntWritable.class);
29
30        job.setMapperClass(Map.class);
31        job.setCombinerClass(Reduce.class);
32    }
33
34 }

```

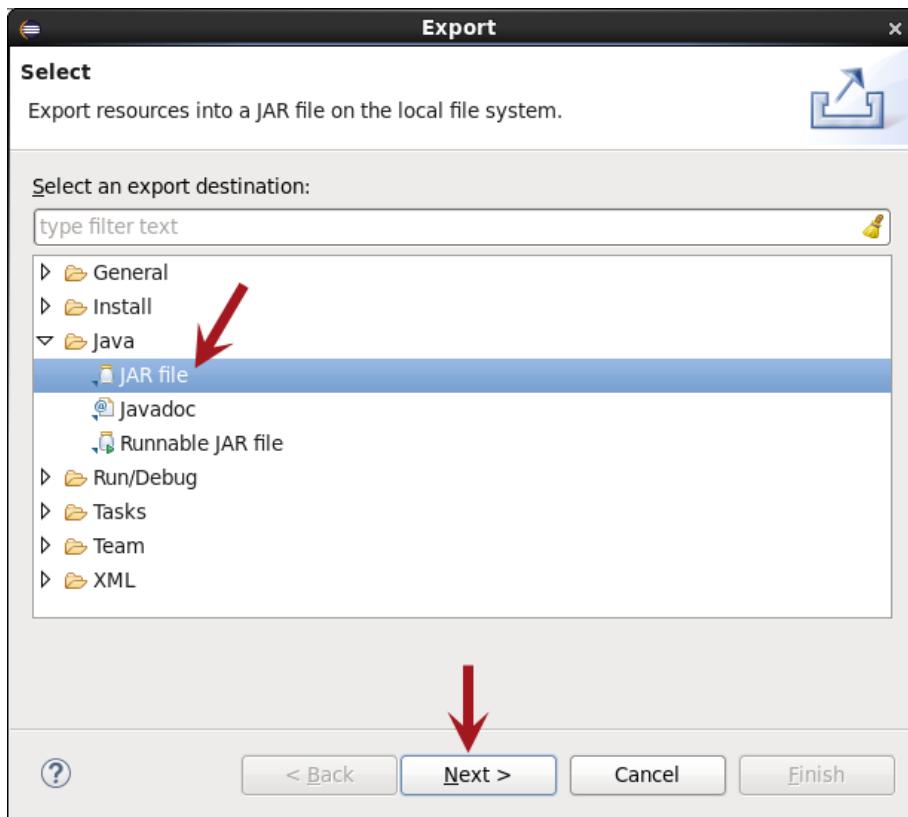


Ensure that there is no compilation error, before exporting the JAR file.

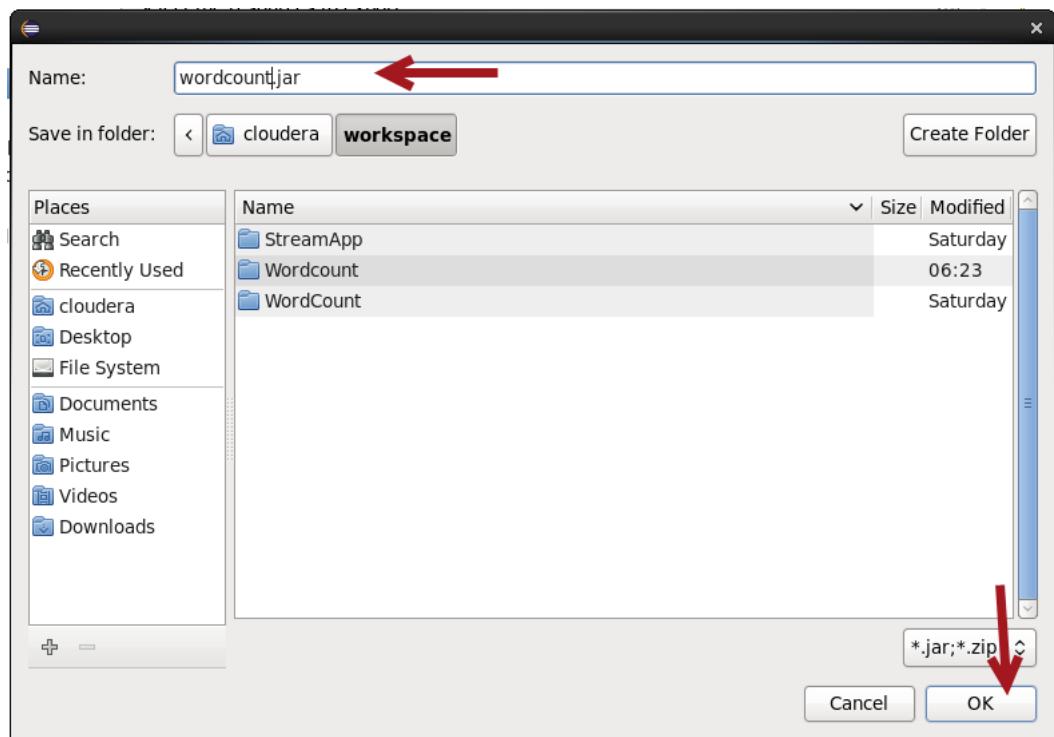
**Step17:** Exporting the JAR file, for that Right Click on WordCount project and select “Export”



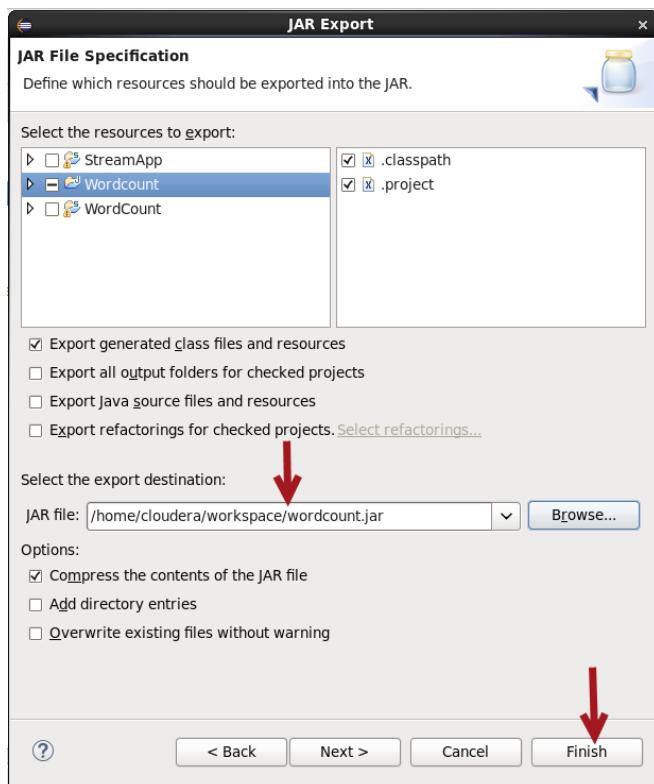
**Step 18:** In the next panel, expand Java and select JAR file



**Step19:** Name the JAR as “WordCount.jar” and select OK.



**Step20:** Browse to the location where you want to save the JAR and select FINISH.



**Step21:** Verify the file at the target location from the terminal.

```
[cloudera@quickstart workspace]$ ls -ltr
total 20
drwxrwxr-x 7 cloudera cloudera 4096 Sep 10 01:43 StreamApp
drwxrwxr-x 8 cloudera cloudera 4096 Sep 10 01:43 WordCount
drwxrwxr-x 5 cloudera cloudera 4096 Sep 13 06:23 Wordcount
-rw-rw-r-- 1 cloudera cloudera 6055 Sep 13 06:40 wordcount.jar
[cloudera@quickstart workspace]$ █
```

➤ Let's run a simple MapReduce WordCount program.

**Step1:** Import the data from local filesystem to HDFS

Under home directory of local filesystem, locate a folder named “datasets” followed by a file named story.txt

```
[cloudera@quickstart ~]$ cd datasets/
[cloudera@quickstart datasets]$ ls -ltr story.txt
-rw-rw-r-- 1 cloudera cloudera 2249 Sep  6 19:28 story.txt
[cloudera@quickstart datasets]$ pwd
/home/cloudera/datasets
[cloudera@quickstart datasets]$ █
```

**Step2:** Create a directory named “data” in the HDFS and import “story.txt” file.

```
| hdfs dfs -mkdir data  
| hdfs dfs -put story.txt data
```

```
[cloudera@quickstart ~]$ hdfs dfs -ls data/story.txt  
-rw-r--r-- 1 cloudera cloudera 2249 2016-09-06 19:29 data/story.txt  
[cloudera@quickstart ~]$ █
```

**Step3:** Now run the program by using the wordCount.jar created in previous steps

```
| hadoop jar wordcount.jar Wordcount data/story.txt output_map
```

```
[cloudera@quickstart workspace]$ pwd  
/home/cloudera/workspace  
[cloudera@quickstart workspace]$ jar -tf wordcount.jar  
META-INF/MANIFEST.MF  
.project  
Wordcount$Map.class  
Wordcount$Reduce.class  
Wordcount.class  
.classpath  
[cloudera@quickstart workspace]$ hadoop jar wordcount.jar Wordcount data/story.txt output_map
```

To check the classes

```
| 16/09/13 06:51:07 INFO mapreduce.Job: map 100% reduce 100%  
16/09/13 06:51:08 INFO mapreduce.Job: Job job_1473496582710_0001 completed  
successfully  
16/09/13 06:51:08 INFO mapreduce.Job: Counters: 49  
File System Counters  
    FILE: Number of bytes read=2906  
    FILE: Number of bytes written=234627  
    FILE: Number of read operations=0  
    FILE: Number of large read operations=0  
    FILE: Number of write operations=0  
    HDFS: Number of bytes read=2374  
    HDFS: Number of bytes written=2035  
    HDFS: Number of read operations=6  
    HDFS: Number of large read operations=0  
    HDFS: Number of write operations=2  
Job Counters  
    Launched map tasks=1  
    Launched reduce tasks=1  
    Data-local map tasks=1  
    Total time spent by all maps in occupied slots (ms)=10457  
    Total time spent by all reduces in occupied slots (ms)=11396  
    Total time spent by all map tasks (ms)=10457  
    Total time spent by all reduce tasks (ms)=11396  
    Total vcore-seconds taken by all map tasks=10457  
    Total vcore-seconds taken by all reduce tasks=11396
```

```
Total megabyte-seconds taken by all map tasks=10707968
Total megabyte-seconds taken by all reduce tasks=11669504
Map-Reduce Framework
  Map input records=17
  Map output records=354
  Map output bytes=3656
  Map output materialized bytes=2906
  Input split bytes=125
  Combine input records=354
  Combine output records=217
  Reduce input groups=217
  Reduce shuffle bytes=2906
  Reduce input records=217
  Reduce output records=217
  Spilled Records=434
  Shuffled Maps =1
  Failed Shuffles=0
  Merged Map outputs=1
  GC time elapsed (ms)=246
  CPU time spent (ms)=1530
  Physical memory (bytes) snapshot=385163264
  Virtual memory (bytes) snapshot=3007664128
  Total committed heap usage (bytes)=354357248
Shuffle Errors
  BAD_ID=0
  CONNECTION=0
  IO_ERROR=0
  WRONG_LENGTH=0
  WRONG_MAP=0
  WRONG_REDUCE=0
File Input Format Counters
  Bytes Read=2249
File Output Format Counters
  Bytes Written=2035
[clooudera@quickstart workspace]$
```

Check the output file:

```
[clooudera@quickstart workspace]$ hdfs dfs -ls output_map
Found 2 items
-rw-r--r--  1 cloudera cloudera      0 2016-09-13 06:51 output_map/_SUCCESS
-rw-r--r--  1 cloudera cloudera 2035 2016-09-13 06:51 output_map/part-r-0
0000
```

The same counters and job status can be checked from Resource Manager URL  
<http://quickstart.cloudera:8088/cluster>

The screenshot shows the 'All Applications' section of the Cloudera Manager interface. The URL in the browser is `quickstart.cloudera:8088/cluster`. The main content area displays 'Cluster Metrics' and 'User Metrics for dr.who'. Below these, a table lists application details. A red arrow points to the URL bar, and another red arrow points to the first row of the application table.

ID	User	Name	Application Type	Queue	StartTime	FinishTime
application_1473496582710_0001	cloudera	Wordcount	MAPREDUCE	root.cloudera	Tue Sep 13 06:50:20 -0700 2016	Tue Sep 13 06:51:06 -0700 2016

## Exercise 4: Launching Spark

- Spark shell can be launched in two ways i.e. Scala and Python.
- To launch Scala spark shell, follow below steps

```
$ spark-shell
```

```
[cloudera@quickstart ~]$ spark-shell
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/lib/zookeeper/lib/slf4j-log4j12-1.7.5.jar!/org/
    slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/jars/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/St
    aticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel).
Welcome to
```

```
    / \
   / \ 
  / \ 
 / \ / \ / \ / \ / \
/ \ / \ . \ \ / \ / \ \ / \
/ \ / \ / \ / \ / \ / \ / \
version 1.6.0
```

```
Using Scala version 2.10.5 (Java HotSpot(TM) 64-Bit Server VM, Java 1.7.0_67)
```



You will see several INFO and WARNING message on the command prompt after launching spark-shell, which can be disregarded.

Logs will appear as below and finally you will get Scala Prompt

SQL context available as sqlContext.

```
scala> 
```

- Spark creates a SparkContext object called sc, verify that the object exists

```
scala> sc
res0: org.apache.spark.SparkContext = org.apache.spark.SparkContext@a23c46d
```

- To know various SparkContext methods which are available, type sc. (sc followed by dot) and then TAB key

---

```
scala> sc.
accumable          accumableCollection
accumulator        addfile
addJar            addSparkListener
appName           applicationAttemptId
applicationId    asInstanceOf
binaryFiles       binaryRecords
broadcast         cancelAllJobs
cancelJobGroup   clearCallSite
clearFiles        clearJars
clearJobGroup    defaultMinPartitions
defaultMinSplits  defaultParallelism
emptyRDD          externalBlockStoreFolderName
files             getAllPools
getCheckpointDir getConf
getExecutorMemoryStatus getExecutorStorageStatus
getLocalProperty  getPersistentRDDs
getPoolForName   getRDDStorageInfo
getSchedulingMode hadoopConfiguration
hadoopFile        hadoopRDD
initLocalProperties isInstanceOf
isLocal           isStopped
jars              killExecutor
killExecutors     makeRDD
```

Spark is based on the concept of Resilient Distributed Dataset (RDD), which is fault tolerant collection of elements that can be operated in parallel.

#### Two ways to create RDDs:

- Parallelizing an existing collection
- Referencing a dataset from an External Storage System

#### Parallelized collection:

To create a parallelized collection holding numbers 1 to 6, use **sc.parallelize** method

```
| val data = Array (1,2,3,4,5,6)
```

```
| val distData = sc.parallelize(data)
```

Once ‘distData’ RDD is created, we can perform **Action** such as:

- Finding the sum, mean & variance

```
| distData.sum
```

```
| distData.mean
```

```
| distData.variance
```

### External Storage System:

Spark can create distributed datasets from any storage system supported by Hadoop, Spark supports text files, Sequence Files and any other Hadoop Input Format.

Load a file from your Local Filesystem say batting.txt using **sc.textFile** method.

```
| val textFile = sc.textFile("file:/home/cloudera/datasets/batting.txt")
```

Operations on RDD will be discussed in next section.

## Exploring Data using RDD operations:

### Transformations

#### **filter**

Return a new dataset formed by selecting those elements of the source on which *func* returns true.

```
| val a = sc.parallelize(1 to 10)
| val b = a.filter(_ % 2 == 0)
| b.collect
```

```

scala> val a = sc.parallelize(1 to 10)
a: org.apache.spark.rdd.RDD[Int] = ParallelCollectionRDD[0] at parallelize at <console>:21

scala> val b = a.filter(_ % 2 == 0)
b: org.apache.spark.rdd.RDD[Int] = MapPartitionsRDD[1] at filter at <console>:23

scala> b.collect
16/08/30 06:59:25 INFO spark.SparkContext: Starting job: collect at <console>:26
16/08/30 06:59:28 INFO scheduler.DAGScheduler: ResultStage 0 (collect at <console>:26) finished in 0.316 s
16/08/30 06:59:28 INFO scheduler.TaskSetManager: Finished task 0.0 in stage 0.0 (TID 0) in 233 ms on localhost (1/1)
16/08/30 06:59:28 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 0.0, whose tasks have all completed, from pool
16/08/30 06:59:28 INFO scheduler.DAGScheduler: Job 0 finished: collect at <console>:26, took 2.285340 s
res0: Array[Int] = Array(2, 4, 6, 8, 10)

```

## map

Return a new distributed dataset formed by passing each element of the source through a function *func*.

```

val a = sc.parallelize(List("dog", "salmon", "salmon", "rat", "elephant"))

val b = a.map(_.length)

val c = a.zip(b)

c.collect

```

```

scala> val a = sc.parallelize(List("dog", "salmon", "salmon", "rat", "elephant"))
a: org.apache.spark.rdd.RDD[String] = ParallelCollectionRDD[2] at parallelize at <console>:21

scala> val b = a.map(_.length)
b: org.apache.spark.rdd.RDD[Int] = MapPartitionsRDD[3] at map at <console>:23

scala> val c = a.zip(b)
c: org.apache.spark.rdd.RDD[(String, Int)] = ZippedPartitionsRDD2[4] at zip at <console>:25

scala> c.collect
16/08/30 07:03:58 INFO spark.SparkContext: Starting job: collect at <console>:28

```

```

16/08/30 07:03:58 INFO executor.Executor: Finished task 0.0 in stage 1.0 (TID 1). 1147 bytes result sent to driver
16/08/30 07:03:58 INFO scheduler.DAGScheduler: ResultStage 1 (collect at <console>:28) finished in 0.062 s
16/08/30 07:03:58 INFO scheduler.DAGScheduler: Job 1 finished: collect at <console>:28, took 0.108580 s
16/08/30 07:03:58 INFO scheduler.TaskSetManager: Finished task 0.0 in stage 1.0 (TID 1) in 64 ms on localhost (1/1)
res1: Array[(String, Int)] = Array((dog,3), (salmon,6), (salmon,6), (rat,3), (elephant,8))

```

## distinct

Return a new dataset that contains the distinct elements of the source dataset.

```

val c = sc.parallelize(List("Gnu", "Cat", "Rat", "Dog", "Gnu", "Rat"), 2)

c.distinct.collect

```

```

scala> val c = sc.parallelize(List("Gnu", "Cat", "Rat", "Dog", "Gnu", "Rat"), 2)
c: org.apache.spark.rdd.RDD[String] = ParallelCollectionRDD[5] at parallelize at <console>:21

scala> c.distinct.collect
16/08/30 07:07:16 INFO spark.SparkContext: Starting job: collect at <console>:24

```

```

16/08/30 07:07:19 INFO executor.Executor: Finished task 1.0 in stage 3.0 (TID 5). 1171 bytes result sent to driver
16/08/30 07:07:19 INFO scheduler.TaskSetManager: Finished task 1.0 in stage 3.0 (TID 5) in 76 ms on localhost (2/2)
16/08/30 07:07:19 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 3.0, whose tasks have all completed, from pool
16/08/30 07:07:19 INFO scheduler.DAGScheduler: ResultStage 3 (collect at <console>:24) finished in 0.385 s
16/08/30 07:07:19 INFO scheduler.DAGScheduler: Job 2 finished: collect at <console>:24, took 2.185199 s
res2: Array[String] = Array(Dog, Cat, Gnu, Rat)

```

## cartesian

When called on datasets of types T and U, returns a dataset of (T, U) pairs (all pairs of elements).

```

-----[{"val x = sc.parallelize(List(1,2,3,4,5))","val y = sc.parallelize(List(6,7,8,9,10))","x.cartesian(y).collect"}]
-----
```

```

scala> val x = sc.parallelize(List(1,2,3,4,5))
x: org.apache.spark.rdd.RDD[Int] = ParallelCollectionRDD[9] at parallelize at <console>:21

scala> val y = sc.parallelize(List(6,7,8,9,10))
y: org.apache.spark.rdd.RDD[Int] = ParallelCollectionRDD[10] at parallelize at <console>:21

scala> x.cartesian(y).collect
16/08/30 07:10:11 INFO spark.SparkContext: Starting job: collect at <console>:26

```

```

16/08/30 07:10:11 INFO scheduler.DAGScheduler: ResultStage 4 (collect at <console>:26) finished in 0.145 s
16/08/30 07:10:11 INFO scheduler.DAGScheduler: Job 3 finished: collect at <console>:26, took 0.264577 s
16/08/30 07:10:11 INFO scheduler.TaskSetManager: Finished task 0.0 in stage 4.0 (TID 6) in 153 ms on localhost (1/1)
16/08/30 07:10:11 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 4.0, whose tasks have all completed, from pool
res3: Array[(Int, Int)] = Array((1,6), (1,7), (1,8), (1,9), (1,10), (2,6), (2,7), (2,8), (2,9), (2,10), (3,6), (3,7), (3,8), (3,9), (3,10), (4,6), (4,7), (4,8), (4,9), (4,10), (5,6), (5,7), (5,8), (5,9), (5,10))

```

## coalesce

Decrease the number of partitions in the RDD to numPartitions. Useful for running operations more efficiently after filtering down a large dataset.

```

-----[{"val y = sc.parallelize(1 to 10, 10)","val z = y.coalesce(2, false)","z.partitions.length"}]
-----
```

```

scala> val y = sc.parallelize(1 to 10, 10)
y: org.apache.spark.rdd.RDD[Int] = ParallelCollectionRDD[12] at parallelize at <console>:21

scala> val z = y.coalesce(2, false)
z: org.apache.spark.rdd.RDD[Int] = CoalescedRDD[13] at coalesce at <console>:23

scala> z.partitions.length
res4: Int = 2

```

## filterByRange

Returns an RDD containing only the items in the key range specified.

```
val randRDD = sc.parallelize(List((2, "cat"), (6, "mouse"), (7, "cup"), (3, "book"), (4, "tv"), (1, "screen"), (5, "heater"))), 3)

val sortedRDD = randRDD.sortByKey()

sortedRDD.filterByRange(1, 3).collect
```

```
16/08/30 07:14:40 INFO executor.Executor: Finished task 0.0 in stage 7.0 (TID 13). 1357 bytes result sent to driver
16/08/30 07:14:40 INFO scheduler.DAGScheduler: ResultStage 7 (collect at <console>:26) finished in 0.080 s
16/08/30 07:14:40 INFO scheduler.DAGScheduler: Job 5 finished: collect at <console>:26, took 0.527731 s
res5: Array[(Int, String)] = Array((1,screen), (2,cat), (3,book))
```

## flatMap

Similar to map, but each input item can be mapped to 0 or more output items (so *func* should return a Seq rather than a single item).

```
val a = sc.parallelize(1 to 10, 5)

a.flatMap(_).collect
```

```
16/08/30 07:18:44 INFO executor.Executor: Finished task 4.0 in stage 9.0 (TID 20). 974 bytes result sent to
16/08/30 07:18:44 INFO scheduler.TaskSetManager: Finished task 4.0 in stage 9.0 (TID 20) in 24 ms on localho
16/08/30 07:18:44 INFO scheduler.DAGScheduler: ResultStage 9 (collect at <console>:24) finished in 0.141 s
16/08/30 07:18:44 INFO scheduler.DAGScheduler: Job 7 finished: collect at <console>:24, took 0.193634 s
res6: Array[Int] = Array(1, 1, 2, 1, 2, 3, 1, 2, 3, 4, 1, 2, 3, 4, 5, 1, 2, 3, 4, 5, 6, 1, 2, 3, 4, 5, 6, 7,
, 8, 1, 2, 3, 4, 5, 6, 7, 8, 9, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10)
```

```
val b = sc.parallelize(List(1, 2, 3), 2).flatMap(x => List(x, x, x)).collect

res85: Array[Int] = Array(1, 1, 1, 2, 2, 2, 3, 3, 3)
```

```
16/08/30 07:17:29 INFO executor.Executor: Running task 1.0 in stage 8.0 (TID 15)
16/08/30 07:17:29 INFO executor.Executor: Finished task 1.0 in stage 8.0 (TID 15). 922 bytes result sent to driver
16/08/30 07:17:29 INFO scheduler.DAGScheduler: ResultStage 8 (collect at <console>:23) finished in 0.104 s
16/08/30 07:17:29 INFO scheduler.DAGScheduler: Job 6 finished: collect at <console>:23, took 0.126929 s
b: Array[Int] = Array(1, 1, 1, 2, 2, 2, 3, 3, 3)
```

```
val lines = sc.parallelize(List("hello world", "hi"))

val words = lines.flatMap(line => line.split(" "))

words.first() // returns "hello"
```

```
16/08/30 07:21:01 INFO scheduler.DAGScheduler: ResultStage 10 (first at <console>:26) finished in 0.006 s
16/08/30 07:21:01 INFO scheduler.DAGScheduler: Job 8 finished: first at <console>:26, took 0.035576 s
res7: String = hello
```

## groupBy

```
| val a = sc.parallelize(1 to 9, 3)  
| a.groupBy(x => { if (x % 2 == 0) "even" else "odd" }).collect
```

```
16/08/30 07:23:14 INFO scheduler.DAGScheduler: ResultStage 12 (collect at <console>:26) finished in 0.182 s  
16/08/30 07:23:14 INFO scheduler.DAGScheduler: Job 9 finished: collect at <console>:26, took 0.447740 s  
res8: Array[(String, Iterable[Int])] = Array((even,CompactBuffer(2, 4, 6, 8)), (odd,CompactBuffer(1, 3, 5, 7,  
9)))
```

## keys

Extracts the keys from all contained tuples and returns them in a new RDD.

```
| val a = sc.parallelize(List("dog", "tiger", "lion", "cat", "panther",  
"eagle"), 2)  
| val b = a.map(x => (x.length, x))  
| b.keys.collect
```

```
16/08/30 07:25:04 INFO scheduler.DAGScheduler: ResultStage 13 (collect at <console>:28) finished in 0.089 s  
16/08/30 07:25:04 INFO scheduler.DAGScheduler: Job 10 finished: collect at <console>:28, took 0.131419 s  
res9: Array[Int] = Array(3, 5, 4, 3, 7, 5)
```

## union

```
| val seta = sc.parallelize(1 to 10)  
| val setb = sc.parallelize(5 to 15)  
| (seta union setb).collect
```

```
16/08/30 07:26:33 INFO scheduler.DAGScheduler: ResultStage 14 (collect at <console>:26) finished in 0.143 s  
16/08/30 07:26:33 INFO scheduler.DAGScheduler: Job 11 finished: collect at <console>:26, took 0.318252 s  
res10: Array[Int] = Array(1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15)
```

## zip

Joins two RDDs by combining the i-th of either partition with each other. The resulting RDD will consist of two-component tuples which are interpreted as key-value pairs by the methods provided by the PairRDDFunctions extension.

```
| val a = sc.parallelize(1 to 100, 3)  
| val b = sc.parallelize(101 to 200, 3)  
| a.zip(b).collect
```

```

16/08/30 07:30:02 INFO scheduler.TaskSetManager: Finished task 2.0 in stage 15.0 (TID 34) in 35 ms on localhost (3/3)
16/08/30 07:30:02 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 15.0, whose tasks have all completed, from pool
res11: Array[(Int, Int)] = Array((1,101), (2,102), (3,103), (4,104), (5,105), (6,106), (7,107), (8,108), (9,109), (10,110), (11,111), (12,112), (13,113), (14,114), (15,115), (16,116), (17,117), (18,118), (19,119), (20,120), (21,121), (22,122), (23,123), (24,124), (25,125), (26,126), (27,127), (28,128), (29,129), (30,130), (31,131), (32,132), (33,133), (34,134), (35,135), (36,136), (37,137), (38,138), (39,139), (40,140), (41,141), (42,142), (43,143), (44,144), (45,145), (46,146), (47,147), (48,148), (49,149), (50,150), (51,151), (52,152), (53,153), (54,154), (55,155), (56,156), (57,157), (58,158), (59,159), (60,160), (61,161), (62,162), (63,163), (64,164), (65,165), (66,166), (67,167), (68,168), (69,169), (70,170), (71,171), (72,172), (73,173), (74,174), (75,175), (76,176), (77,177), (78...
scala> 
```

## Action

### collect

Return all the elements of the dataset as an array at the driver program. This is usually useful after a filter or other operation that returns a sufficiently small subset of the data.

```

val c = sc.parallelize(List("Gnu", "Cat", "Rat", "Dog", "Gnu", "Rat"), 2)
c.collect

```

```

16/08/30 07:31:52 INFO scheduler.TaskSetManager: Finished task 1.0 in stage 16.0 (TID 36) in 7 ms on localhost (2/2)
16/08/30 07:31:52 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 16.0, whose tasks have all completed, from pool
16/08/30 07:31:52 INFO scheduler.DAGScheduler: ResultStage 16 (collect at <console>:24) finished in 0.010 s
16/08/30 07:31:52 INFO scheduler.DAGScheduler: Job 13 finished: collect at <console>:24, took 0.018261 s
res12: Array[String] = Array(Gnu, Cat, Rat, Dog, Gnu, Rat)

```

### collectAsMap

Similar to collect, but works on key-value RDDs and converts them into Scala maps to preserve their key-value structure.

```

val a = sc.parallelize(List(1, 2, 1, 3), 1)
val b = a.zip(a)
b.collectAsMap

```

```

16/08/30 07:33:33 INFO scheduler.TaskSetManager: Finished task 0.0 in stage 17.0 (TID 37) in 145 ms on localhost (1/1)
16/08/30 07:33:33 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 17.0, whose tasks have all completed, from pool
res13: scala.collection.Map[Int,Int] = Map(2 -> 2, 1 -> 1, 3 -> 3)

```

## count

Return the number of elements in the dataset.

```
| val c = sc.parallelize(List("Gnu", "Cat", "Rat", "Dog"), 2)
| c.count
| res2: Long = 4
```

```
16/08/30 07:34:52 INFO scheduler.TaskSetManager: Finished task 1.0 in stage 18.0 (TID 39) in 11 ms on localhost (2/2)
16/08/30 07:34:52 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 18.0, whose tasks have all completed, from pool
16/08/30 07:34:52 INFO scheduler.DAGScheduler: ResultStage 18 (count at <console>:24) finished in 0.040 s
16/08/30 07:34:52 INFO scheduler.DAGScheduler: Job 15 finished: count at <console>:24, took 0.054256 s
res14: Long = 4
```

## reduce

Aggregate the elements of the dataset using a function *func* (which takes two arguments and returns one). The function should be commutative and associative so that it can be computed correctly in parallel.

```
| val a = sc.parallelize(1 to 100, 3)
| a.reduce(_+_)
```

```
16/08/30 08:32:49 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 19.0, whose tasks have all completed, from pool
16/08/30 08:32:49 INFO scheduler.DAGScheduler: ResultStage 19 (reduce at <console>:24) finished in 0.099 s
16/08/30 08:32:49 INFO scheduler.DAGScheduler: Job 16 finished: reduce at <console>:24, took 0.315500 s
res15: Int = 5050
```

## take

Return an array with the first *n* elements of the dataset.

```
| val b = sc.parallelize(List("dog", "cat", "ape", "salmon", "gnu"), 2)
| b.take(2)
```

```
16/08/30 08:34:15 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 20.0, whose tasks have all completed, from pool
16/08/30 08:34:15 INFO scheduler.DAGScheduler: ResultStage 20 (take at <console>:24) finished in 0.035 s
16/08/30 08:34:15 INFO scheduler.DAGScheduler: Job 17 finished: take at <console>:24, took 0.142047 s
res16: Array[String] = Array(dog, cat)
```

```
| val b = sc.parallelize(1 to 100, 5)
| b.take(30)
```

```
16/08/30 08:35:57 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 22.0, whose tasks have all completed, from pool
16/08/30 08:35:58 INFO scheduler.DAGScheduler: ResultStage 22 (take at <console>:24) finished in 0.004 s
16/08/30 08:35:58 INFO scheduler.DAGScheduler: Job 19 finished: take at <console>:24, took 0.064874 s
res17: Array[Int] = Array(1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23,
24, 25, 26, 27, 28, 29, 30)
```

## first

Return the first element of the dataset (similar to take(1)).

```
val c = sc.parallelize(List("Gnu", "Cat", "Rat", "Dog"), 2)
c.first
```

```
16/08/30 08:37:15 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 23.0, whose tasks have all completed, from pool
16/08/30 08:37:15 INFO scheduler.DAGScheduler: ResultStage 23 (first at <console>:24) finished in 0.006 s
16/08/30 08:37:15 INFO scheduler.DAGScheduler: Job 20 finished: first at <console>:24, took 0.016040 s
res18: String = Gnu
```

## countByValue

Returns a map that contains all unique values of the RDD and their respective occurrence counts.

```
val b = sc.parallelize(List(1,2,3,4,5,6,7,8,2,4,2,1,1,1,1,1))
b.countByValue
```

```
16/08/30 08:38:18 INFO scheduler.DAGScheduler: ResultStage 25 (countByValue at <console>:24) finished in 0.040 s
16/08/30 08:38:18 INFO scheduler.DAGScheduler: Job 21 finished: countByValue at <console>:24, took 0.495353 s
res19: scala.collection.Map[Int,Long] = Map(5 -> 1, 1 -> 6, 6 -> 1, 2 -> 3, 7 -> 1, 3 -> 1, 8 -> 1, 4 -> 2)
```

## lookup

Scans the RDD for all keys that match the provided value and returns their values as a Scala sequence.

```
val a = sc.parallelize(List("dog", "tiger", "lion", "cat", "panther",
"eagle"), 2)
val b = a.map(x => (x.length, x))
b.lookup(5)
```

```
16/08/30 08:39:47 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 26.0, whose tasks have all completed, from pool  
16/08/30 08:39:47 INFO scheduler.DAGScheduler: ResultStage 26 (lookup at <console>:28) finished in 0.087 s  
16/08/30 08:39:47 INFO scheduler.DAGScheduler: Job 22 finished: lookup at <console>:28, took 0.188596 s  
res20: Seq[String] = WrappedArray(tiger, eagle)
```

## max

Returns the largest element in the RDD

```
val y = sc.parallelize(10 to 30)  
y.max
```

```
16/08/30 08:41:08 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 27.0, whose tasks have all completed, from pool  
16/08/30 08:41:08 INFO scheduler.DAGScheduler: ResultStage 27 (max at <console>:24) finished in 0.042 s  
16/08/30 08:41:08 INFO scheduler.DAGScheduler: Job 23 finished: max at <console>:24, took 0.115601 s  
res21: Int = 30
```

```
val a = sc.parallelize(List((10, "dog"), (3, "tiger"), (9, "lion"), (18, "cat")))  
a.max
```

```
16/08/30 08:42:29 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 28.0, whose tasks have all completed, from pool  
16/08/30 08:42:29 INFO scheduler.DAGScheduler: ResultStage 28 (max at <console>:24) finished in 0.031 s  
16/08/30 08:42:29 INFO scheduler.DAGScheduler: Job 24 finished: max at <console>:24, took 0.054414 s  
res22: (Int, String) = (18,cat)
```

## min

Returns the smallest element in the RDD

```
val y = sc.parallelize(10 to 30)  
y.min
```

```
16/08/30 08:43:49 INFO executor.Executor: Finished task 0.0 in stage 29.0 (TID 53). 1031 bytes result sent to driver  
16/08/30 08:43:49 INFO scheduler.DAGScheduler: ResultStage 29 (min at <console>:24) finished in 0.075 s  
16/08/30 08:43:49 INFO scheduler.DAGScheduler: Job 25 finished: min at <console>:24, took 0.170125 s  
res23: Int = 10
```

```
val a = sc.parallelize(List((10, "dog"), (3, "tiger"), (9, "lion"), (8, "cat")))  
a.min
```

```
16/08/30 08:45:21 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 30.0, whose tasks have all completed, from pool
16/08/30 08:45:21 INFO scheduler.DAGScheduler: ResultStage 30 (min at <console>:24) finished in 0.012 s
16/08/30 08:45:21 INFO scheduler.DAGScheduler: Job 26 finished: min at <console>:24, took 0.020695 s
res24: (Int, String) = (3,tiger)
```

## mean

Calls stats and extracts the mean component.

```
val a = sc.parallelize(List(9.1, 1.0, 1.2, 2.1, 1.3, 5.0, 2.0, 2.1, 7.4,
7.5, 7.6, 8.8, 10.0, 8.9, 5.5), 3)
a.mean
```

```
16/08/30 08:46:49 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 31.0, whose tasks have all completed, from pool
16/08/30 08:46:49 INFO scheduler.DAGScheduler: ResultStage 31 (mean at <console>:24) finished in 0.425 s
16/08/30 08:46:49 INFO scheduler.DAGScheduler: Job 27 finished: mean at <console>:24, took 0.526866 s
res25: Double = 5.3
```

## variance

Calls stats and extracts either variance-component or corrected sampleVariance-component.

```
val a = sc.parallelize(List(9.1, 1.0, 1.2, 2.1, 1.3, 5.0, 2.0, 2.1, 7.4,
7.5, 7.6, 8.8, 10.0, 8.9, 5.5), 3)
a.variance
```

```
16/08/30 08:48:05 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 32.0, whose tasks have all completed, from pool
16/08/30 08:48:05 INFO scheduler.DAGScheduler: ResultStage 32 (variance at <console>:24) finished in 0.018 s
16/08/30 08:48:05 INFO scheduler.DAGScheduler: Job 28 finished: variance at <console>:24, took 0.032808 s
res26: Double = 10.60533333333332
```

# PairRDD

## countByKey

Only available on RDDs of type (K, V). Returns a hashmap of (K, Int) pairs with the count of each key.

```
val c = sc.parallelize(List((3, "Gnu"), (3, "Yak"), (5, "Mouse"), (3, "Dog")), 2)
c.countByKey
```

```

16/08/30 08:49:39 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 34.0, whose tasks have all completed, from pool
16/08/30 08:49:39 INFO scheduler.DAGScheduler: ResultStage 34 (countByKey at <console>:24) finished in 0.073 s
16/08/30 08:49:39 INFO scheduler.DAGScheduler: Job 29 finished: countByKey at <console>:24, took 0.384196 s
res27: scala.collection.Map[Int,Long] = Map(3 -> 3, 5 -> 1)

```

## groupByKey

When called on a dataset of (K, V) pairs, returns a dataset of (K, Iterable<V>) pairs.

```

val a = sc.parallelize(List("dog", "tiger", "lion", "cat", "spider", "eagle"), 2)
val b = a.keyBy(_.length)
b.groupByKey.collect

```

```

16/08/30 08:51:02 INFO executor.Executor: Finished task 1.0 in stage 36.0 (TID 68). 1703 bytes result sent to driver
16/08/30 08:51:02 INFO scheduler.DAGScheduler: ResultStage 36 (collect at <console>:26) finished in 0.042 s
16/08/30 08:51:02 INFO scheduler.DAGScheduler: Job 30 finished: collect at <console>:26, took 0.415007 s
res28: Array[(Int, Iterable[String])] = Array((4,CompactBuffer(lion)), (6,CompactBuffer(spider)), (3,CompactBuffer(dog, cat)), (5,CompactBuffer(tiger, eagle)))

```

## reduceByKey

When called on a dataset of (K, V) pairs, returns a dataset of (K, V) pairs where the values for each key are aggregated using the given reduce function *func*, which must be of type (V,V) => V

```

val a = sc.parallelize(List("dog", "cat", "owl", "gnu", "ant"), 2)
val b = a.map(x => (x.length, x))
b.reduceByKey(_ + _).collect

```

```

16/08/30 08:52:58 INFO scheduler.TaskSetManager: Finished task 1.0 in stage 38.0 (TID 72) in 25 ms on localhost (2/2)
16/08/30 08:52:58 INFO scheduler.DAGScheduler: ResultStage 38 (collect at <console>:28) finished in 0.033 s
16/08/30 08:52:58 INFO scheduler.DAGScheduler: Job 31 finished: collect at <console>:28, took 0.413629 s
res29: Array[(Int, String)] = Array((3,dogcatowlgnuant))

```

## foldByKey

Very similar to *fold*, but performs the folding separately for each key of the RDD. This function is only available if the RDD consists of two-component tuples.

```

val deptEmployees =
List(
  ("dept1", ("kumar1", 1000.0)),

```

```

        ("dept1", ("kumar2",1200.0)),
        ("dept2", ("kumar3",2200.0)),
        ("dept2", ("kumar4",1400.0)),
        ("dept2", ("kumar5",1000.0)),
        ("dept2", ("kumar6",800.0)),
        ("dept1", ("kumar7",2000.0)),
        ("dept1", ("kumar8",1000.0)),
        ("dept1", ("kumar9",500.0))
    )

val employeeRDD = sc.makeRDD(deptEmployees)

val maxByDept =
employeeRDD.foldByKey(("dummy",Double.MinValue))((acc,element) => if(acc._2 > element._2)acc else element)

println("Maximum salaries in each dept" + maxByDept.collect().toList)

```

```

16/08/30 08:56:35 INFO executor.Executor: Finished task 0.0 in stage 40.0 (TID 74). 1375 bytes result sent to driver
16/08/30 08:56:35 INFO scheduler.DAGScheduler: ResultStage 40 (collect at <console>:28) finished in 0.087 s
16/08/30 08:56:35 INFO scheduler.DAGScheduler: Job 32 finished: collect at <console>:28, took 0.351474 s
Maximum salaries in each deptList((dept2,(kumar3,2200.0)), (dept1,(kumar7,2000.0)))

```

## cogroup

When called on datasets of type (K, V) and (K, W), returns a dataset of (K, (Iterable<V>, Iterable<W>)) tuples.

```

val a = sc.parallelize(List(1, 2, 1, 3), 1)

val b = a.map((_, "b"))

val c = a.map((_, "c"))

b.cogroup(c).collect

```

```

16/08/30 08:58:50 INFO executor.Executor: Finished task 0.0 in stage 43.0 (TID 77). 2177 bytes result sent to driver
16/08/30 08:58:50 INFO scheduler.DAGScheduler: ResultStage 43 (collect at <console>:28) finished in 0.249 s
16/08/30 08:58:50 INFO scheduler.DAGScheduler: Job 33 finished: collect at <console>:28, took 0.668327 s
res31: Array[(Int, (Iterable[String], Iterable[String]))] = Array((1,(CompactBuffer(b, b),CompactBuffer(c, c))), (3,(CompactBuffer(b),CompactBuffer(c))), (2,(CompactBuffer(b),CompactBuffer(c))))

```

```

val x = sc.parallelize(List((1, "apple"), (2, "banana"), (3, "orange"), (4, "kiwi")), 2)

val y = sc.parallelize(List((5, "computer"), (1, "laptop"), (1, "desktop"), (4, "iPad")), 2)

x.cogroup(y).collect

```

```

16/08/30 09:00:55 INFO executor.Executor: Finished task 1.0 in stage 46.0 (TID 83). 2189 bytes result sent to driver
16/08/30 09:00:55 INFO scheduler.DAGScheduler: ResultStage 46 (collect at <console>:26) finished in 0.123 s
16/08/30 09:00:55 INFO scheduler.DAGScheduler: Job 34 finished: collect at <console>:26, took 0.624333 s
res32: Array[Int, (Iterable[String], Iterable[String])] = Array((4,(CompactBuffer(kiwi),CompactBuffer(iPad))), (2,(CompactBuffer(banana),CompactBuffer())), (1,(CompactBuffer(apple),CompactBuffer(laptop, desktop))), (3,(CompactBuffer(orange),CompactBuffer())), (5,(CompactBuffer(),CompactBuffer(computer))))

```

## Join

Performs an inner join using two key-value RDDs.

```

val a = sc.parallelize(List("dog", "salmon", "salmon", "rat", "elephant"), 3)
val b = a.keyBy(_.length)
val c =
sc.parallelize(List("dog", "cat", "gnu", "salmon", "rabbit", "turkey", "wolf", "bear", "bee"), 3)
val d = c.keyBy(_.length)
b.join(d).collect

```

```

16/08/30 09:02:41 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 49.0, whose tasks have all completed, from pool
16/08/30 09:02:41 INFO scheduler.DAGScheduler: ResultStage 49 (collect at <console>:30) finished in 0.470 s
16/08/30 09:02:41 INFO scheduler.DAGScheduler: Job 35 finished: collect at <console>:30, took 1.019005 s
res33: Array[Int, (String, String)] = Array((6,(salmon,salmon)), (6,(salmon,rabbit)), (6,(salmon,turkey)), (6,(salmon,salmon)), (6,(salmon,rabbit)), (6,(salmon,turkey)), (3,(dog,dog)), (3,(dog,cat)), (3,(dog,gnu)), (3,(dog,bee)), (3,(rat,dog)), (3,(rat,cat)), (3,(rat,gnu)), (3,(rat,bee)))

```

## leftOuterJoin

Performs an left outer join using two key-value RDDs.

```

val a = sc.parallelize(List("dog", "salmon", "salmon", "rat", "elephant"), 3)
val b = a.keyBy(_.length)
val c =
sc.parallelize(List("dog", "cat", "gnu", "salmon", "rabbit", "turkey", "wolf", "bear", "bee"), 3)
val d = c.keyBy(_.length)
b.leftOuterJoin(d).collect

```

```

16/08/30 09:03:59 INFO scheduler.DAGScheduler: Job 36 finished: collect at <console>:30, took 0.675634 s
16/08/30 09:03:59 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 52.0, whose tasks have all completed, from pool
res34: Array[Int, (String, Option[String])] = Array((6,(salmon,Some(salmon))), (6,(salmon,Some(rabbit))), (6,(salmon,Some(turkey))), (6,(salmon,Some(salmon))), (6,(salmon,Some(rabbit))), (6,(salmon,Some(turkey))), (3,(dog,Some(dog))), (3,(dog,Some(cat))), (3,(dog,Some(gnu))), (3,(dog,Some(bee))), (3,(rat,Some(dog))), (3,(rat,Some(cat))), (3,(rat,Some(gnu))), (3,(rat,Some(bee))), (8,(elephant,None)))

```

## rightOuterJoin

Performs an right outer join using two key-value RDDs.

```
val a = sc.parallelize(List("dog", "salmon", "salmon", "rat", "elephant"), 3)
val b = a.keyBy(_.length)
val c =
sc.parallelize(List("dog", "cat", "gnu", "salmon", "rabbit", "turkey", "wolf", "bear", "bee"), 3)
val d = c.keyBy(_.length)
b.rightOuterJoin(d).collect
```

```
16/08/30 09:05:30 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 55.0, whose tasks have all completed, from pool
16/08/30 09:05:30 INFO scheduler.DAGScheduler: ResultStage 55 (collect at <console>:30) finished in 0.181 s
16/08/30 09:05:30 INFO scheduler.DAGScheduler: Job 37 finished: collect at <console>:30, took 1.003812 s
res35: Array[(Int, (Option[String], String))] = Array((6,(Some(salmon),salmon)), (6,(Some(salmon),rabbit)),
(6,(Some(salmon),turkey)), (6,(Some(salmon),salmon)), (6,(Some(salmon),rabbit)), (6,(Some(salmon),turkey)),
(3,(Some(dog),dog)), (3,(Some(dog),cat)), (3,(Some(dog),gnu)), (3,(Some(dog),bee)), (3,(Some(rat),dog)), (3,
(Some(rat),cat)), (3,(Some(rat),gnu)), (3,(Some(rat),bee)), (4,(None,wolf)), (4,(None,bear)))
```

## keyBy

Constructs two-component tuples (key-value pairs) by applying a function on each data item. The result of the function becomes the key and the original data item becomes the value of the newly created tuples.

```
val a = sc.parallelize(List("dog", "salmon", "salmon", "rat", "elephant"), 3)
val b = a.keyBy(_.length)
b.collect
```

```
16/08/30 09:14:29 INFO scheduler.TaskSetManager: Finished task 2.0 in stage 56.0 (TID 113) in 7 ms on localhost (3/3)
16/08/30 09:14:29 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 56.0, whose tasks have all completed, from pool
16/08/30 09:14:29 INFO scheduler.DAGScheduler: ResultStage 56 (collect at <console>:26) finished in 0.053 s
16/08/30 09:14:29 INFO scheduler.DAGScheduler: Job 38 finished: collect at <console>:26, took 0.073531 s
res36: Array[(Int, String)] = Array((3,dog), (6,salmon), (6,salmon), (3,rat), (8,elephant))
```

## mapValues

Takes the values of a RDD that consists of two-component tuples, and applies the provided function to transform each value. Then, it forms new two-component tuples using the key and the transformed value and stores them in a new RDD.

```
-----  
val a = sc.parallelize(List("dog", "tiger", "lion", "cat", "panther",  
"eagle"), 2)  
  
val b = a.map(x => (x.length, x))  
  
b.mapValues("x" + _ + "x").collect  
-----  
  
16/08/30 09:21:21 INFO executor.Executor: Finished task 1.0 in stage 57.0 (TID 115). 1118 bytes result sent  
to driver  
16/08/30 09:21:21 INFO scheduler.DAGScheduler: ResultStage 57 (collect at <console>:28) finished in 0.082 s  
16/08/30 09:21:21 INFO scheduler.DAGScheduler: Job 39 finished: collect at <console>:28, took 0.182411 s  
res37: Array[(Int, String)] = Array((3,xdogx), (5,xtigerx), (4,xlionx), (3,xcatx), (7,xpantherx), (5,xeaglex))  
-----
```

## cache

The cache( ) method is a shorthand for using the default storage level, which is StorageLevel.MEMORY\_ONLY (store serialized objects in memory)

```
-----  
val c = sc.parallelize(List("Gnu", "Cat", "Rat", "Dog", "Gnu", "Rat"), 2)  
  
c.getStorageLevel  
-----  
  
scala> val c = sc.parallelize(List("Gnu", "Cat", "Rat", "Dog", "Gnu", "Rat"), 2)  
c: org.apache.spark.rdd.RDD[String] = ParallelCollectionRDD[99] at parallelize at <console>:21  
  
scala> c.getStorageLevel  
res38: org.apache.spark.storage.StorageLevel = StorageLevel(false, false, false, false, 1)  
  
-----  
c.cache  
  
c.getStorageLevel  
-----  
  
scala> c.cache  
res39: c.type = ParallelCollectionRDD[99] at parallelize at <console>:21  
  
scala> c.getStorageLevel  
res40: org.apache.spark.storage.StorageLevel = StorageLevel(false, true, false, true, 1)  
-----
```

## repartition

Reshuffle the data in the RDD randomly to create either more or fewer partitions and balance it across them. This always shuffles all data over the network.

```
val rdd = sc.parallelize(List(1, 2, 10, 4, 5, 2, 1, 1, 1), 3)
rdd.partitions.length
val rdd2 = rdd.repartition(5)
rdd2.partitions.length
```

```
scala> val rdd = sc.parallelize(List(1, 2, 10, 4, 5, 2, 1, 1, 1), 3)
rdd: org.apache.spark.rdd.RDD[Int] = ParallelCollectionRDD[102] at parallelize at <console>:21
scala> rdd.partitions.length
res43: Int = 3
scala> val rdd2 = rdd.repartition(5)
rdd2: org.apache.spark.rdd.RDD[Int] = MapPartitionsRDD[106] at repartition at <console>:23
scala> rdd2.partitions.length
res44: Int = 5
```

## Developing with Spark

## REPL

**Example:** Counting the occurrence of lines having a particular word in it

```
val textFile = sc.textFile("file:/home/cloudera/Datasets/Joyce.txt")
```

```
textFile.count() //Return the number of elements in the dataset
```

```
16/08/30 09:34:43 INFO executor.Executor: Finished task 0.0 in stage 59.0 (TID 136). 2082 bytes result sent
to driver
16/08/30 09:34:43 INFO scheduler.DAGScheduler: ResultStage 59 (count at <console>:24) finished in 2.238 s
16/08/30 09:34:43 INFO scheduler.DAGScheduler: Job 41 finished: count at <console>:24, took 2.310447 s
res45: Long = 33056
```

```
textFile.first() //Return the first element of the dataset
```

```
16/08/30 09:35:03 INFO executor.Executor: Finished task 0.0 in stage 60.0 (TID 137). 2101 bytes result sent  
to driver  
16/08/30 09:35:03 INFO scheduler.DAGScheduler: ResultStage 60 (first at <console>:24) finished in 0.035 s  
16/08/30 09:35:03 INFO scheduler.DAGScheduler: Job 42 finished: first at <console>:24, took 0.091601 s  
res46: String = The Project Gutenberg EBook of Ulysses, by James Joyce
```

```
val linesWithJoyce = textFile.filter(line => line.contains("Joyce"))  
linesWithJoyce.count() //How many lines contains "Joyce"
```

```
16/08/30 09:35:26 INFO executor.Executor: Finished task 0.0 in stage 61.0 (TID 138). 2082 bytes result sent  
to driver  
16/08/30 09:35:26 INFO scheduler.DAGScheduler: ResultStage 61 (count at <console>:26) finished in 0.064 s  
16/08/30 09:35:26 INFO scheduler.DAGScheduler: Job 43 finished: count at <console>:26, took 0.094628 s  
res47: Long = 4
```

**Example:** Add up the sizes of all the lines

```
val lineLength = textFile.map(s => s.length)  
val totalLength = lineLength.reduce((a, b) => a + b) //Aggregate the  
elements of the dataset using a function
```

```
16/08/30 09:40:54 INFO executor.Executor: Finished task 0.0 in stage 62.0 (TID 139). 2160 bytes result sent  
to driver  
16/08/30 09:40:54 INFO scheduler.DAGScheduler: ResultStage 62 (reduce at <console>:31) finished in 0.262 s  
16/08/30 09:40:54 INFO scheduler.DAGScheduler: Job 44 finished: reduce at <console>:31, took 0.337723 s  
totalLength: Int = 1506967
```

**Example:** To find out the line with maximum words

```
textFile.map(line => line.split(" ").size).reduce((a,b) => if(a>b) a else  
b)
```

```
16/08/30 09:41:48 INFO executor.Executor: Finished task 0.0 in stage 63.0 (TID 140). 2160 bytes result sent  
to driver  
16/08/30 09:41:48 INFO scheduler.DAGScheduler: ResultStage 63 (reduce at <console>:30) finished in 1.048 s  
16/08/30 09:41:48 INFO scheduler.DAGScheduler: Job 45 finished: reduce at <console>:30, took 1.070554 s  
res48: Int = 22
```

## Zeppelin

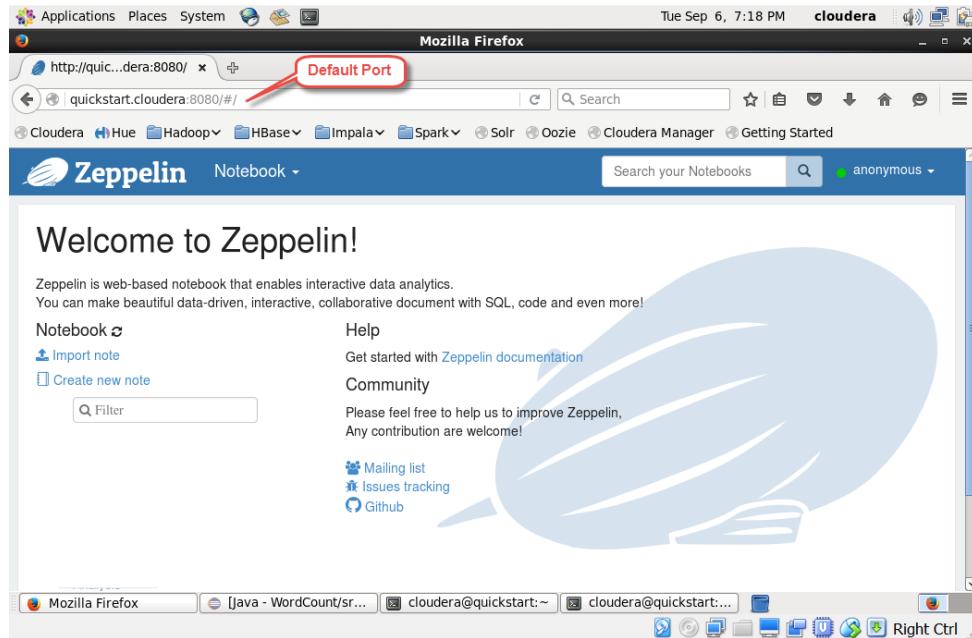
To start Zeppelin service

```
[ - - - - - ] cd /usr/lib/zeppelin-0.6.0-bin-all  
[ - - - - - ] sudo bin/zeppelin-daemon.sh start
```

```
[cloudera@quickstart ~]$ cd /usr/lib/zeppelin-0.6.0-bin-all/  
[cloudera@quickstart zeppelin-0.6.0-bin-all]$ ls bin  
common.cmd      functions.sh          interpreter.sh      zeppelin.sh  
common.sh       install-interpreter.sh zeppelin.cmd  
functions.cmd   interpreter.cmd      zeppelin-daemon.sh  
[cloudera@quickstart zeppelin-0.6.0-bin-all]$ bin/zeppelin-daemon.sh start  
Zeppelin start                                         [ OK ]  
[cloudera@quickstart zeppelin-0.6.0-bin-all]$ █
```

## Accessing Zeppelin UI

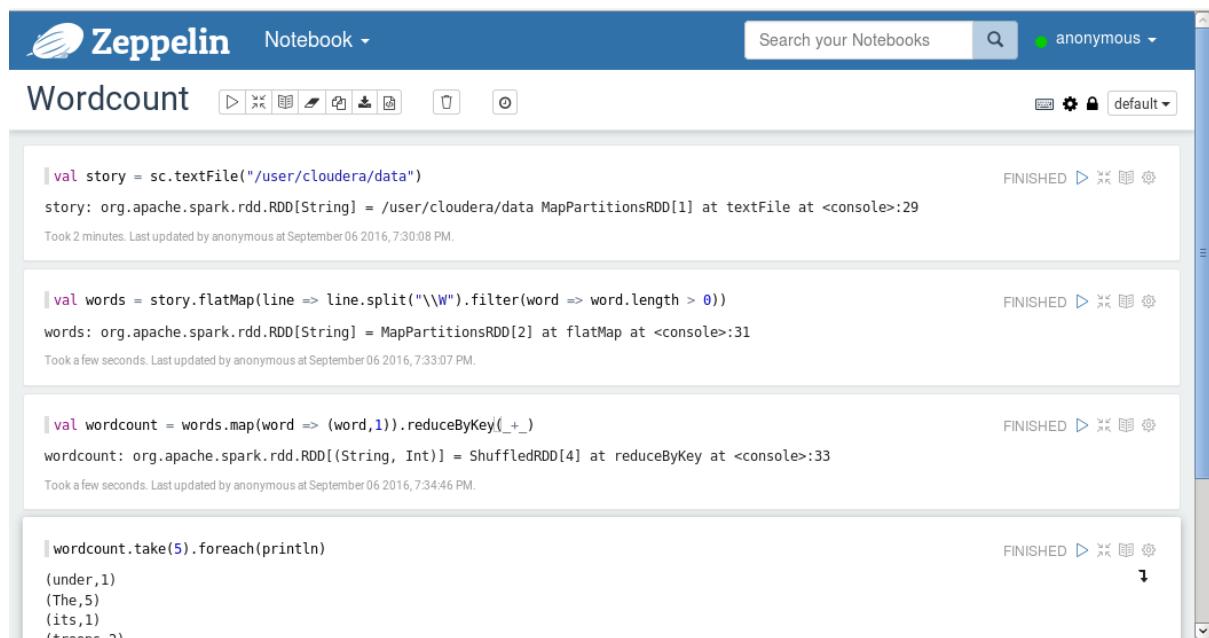
Open browser and use port 8080 to access Zeppelin notebook <quickstart.cloudera:8080>



## WordCount Example

Create a new note, name it as ‘WordCount’ and run following commands from zeppelin notebook.

```
val story = sc.textFile("/user/cloudera/data")  
  
val words = story.flatMap(line => line.split("\\W").filter(word =>  
word.length > 0))  
  
val wordcount = words.map(word => (word,1)).reduceByKey(_+_)  
  
wordcount.take(5).foreach(println)
```



The screenshot shows the Zeppelin Notebook interface with the title "Wordcount". The execution history consists of four cells:

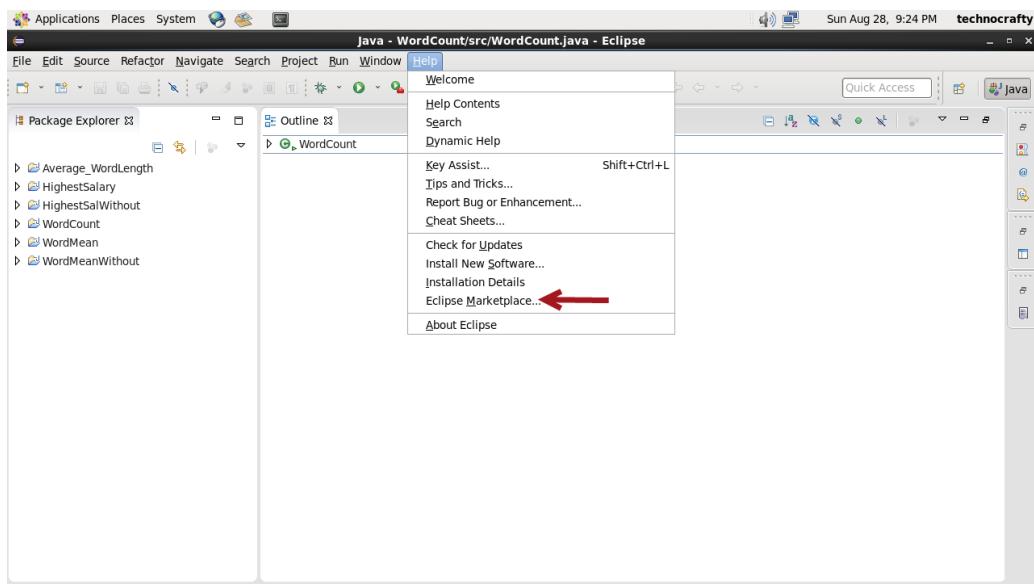
- Cell 1:** `val story = sc.textFile("/user/cloudera/data")` - FINISHED
- Cell 2:** `val words = story.flatMap(line => line.split("\\W").filter(word => word.length > 0))` - FINISHED
- Cell 3:** `val wordcount = words.map(word => (word,1)).reduceByKey(_+_)` - FINISHED
- Cell 4:** `wordcount.take(5).foreach(println)` - FINISHED, output: (under,1), (The,5), (its,1), (+more...)

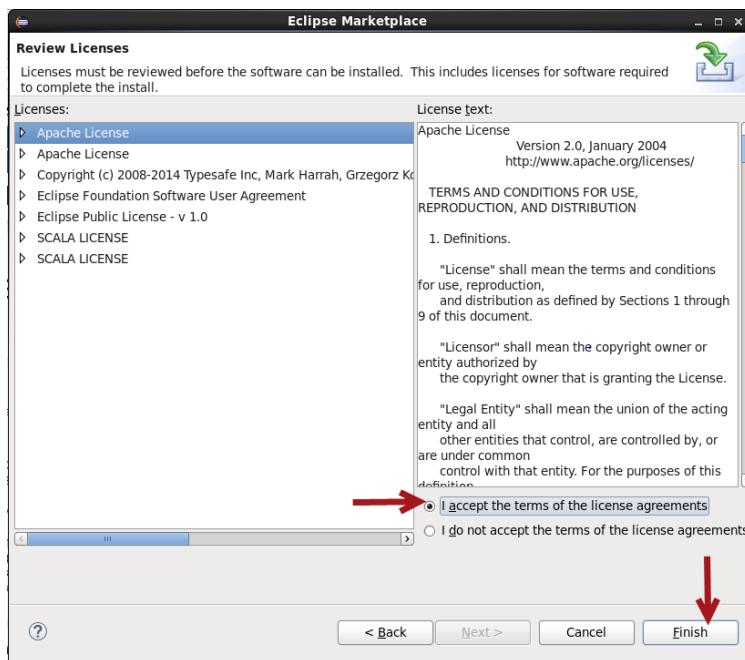
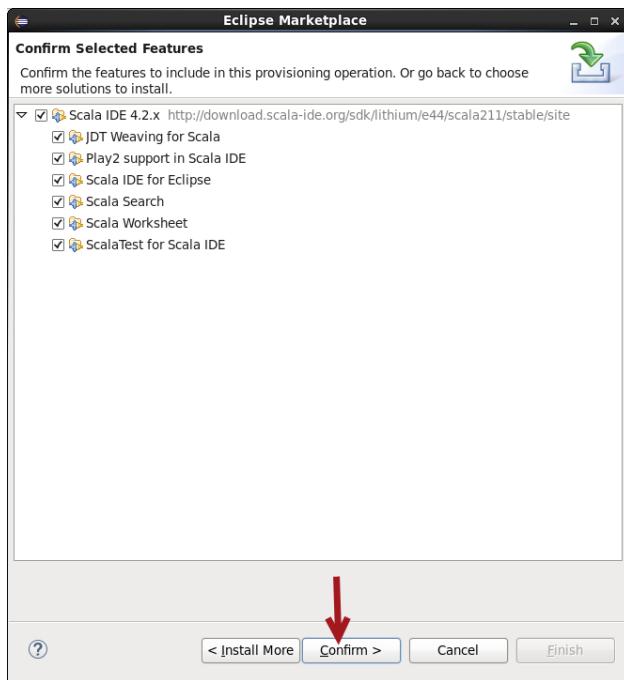
## Exercise 5: Building Spark Application using Eclipse IDE

**Goal:** Create a scala project for Spark. Use maven for lib management.

1. Install Scala IDE for Eclipse from Eclipse Marketplace

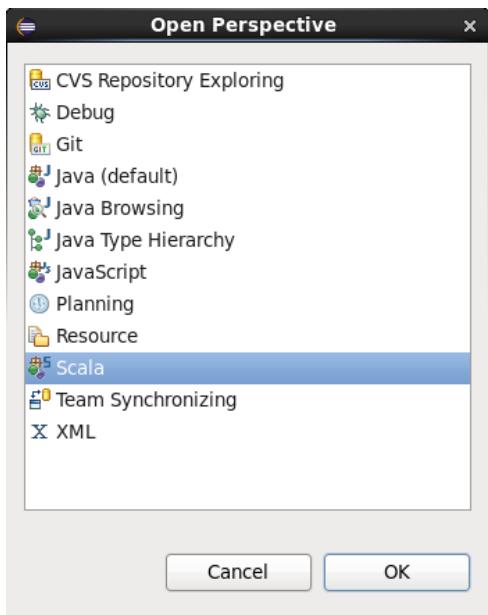
Note: Scala IDE is already installed for this training.



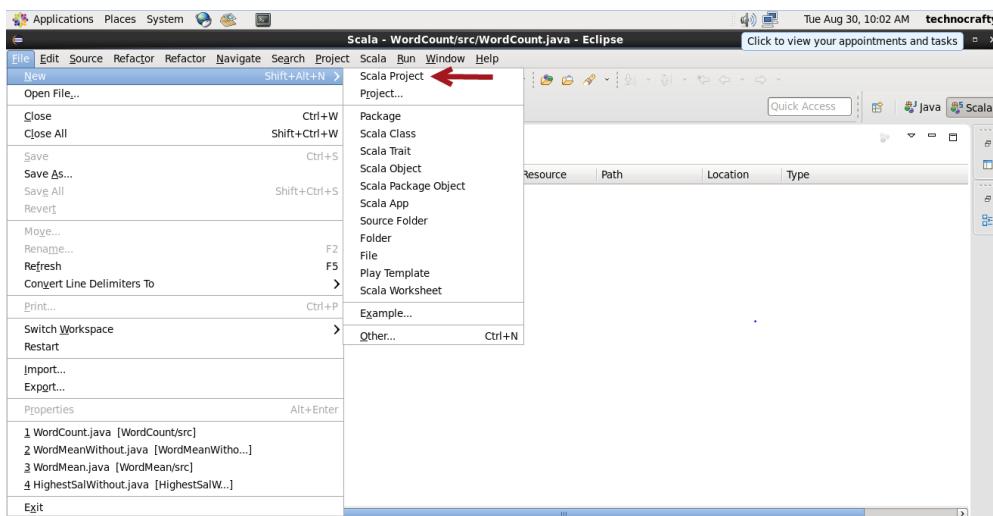


## 2. Switch to Scala perspective.

Window > Perspectives > Open Perspective > Other ... > Scala > Click OK



3. Create new Scala Project using the menu bar option File > New > Scala Project. Give a name of your choice for the project.



4. Right click on the project and configure it as a Maven project. No need change any configuration, however group Id usually takes form like com.mycompany.projectname.

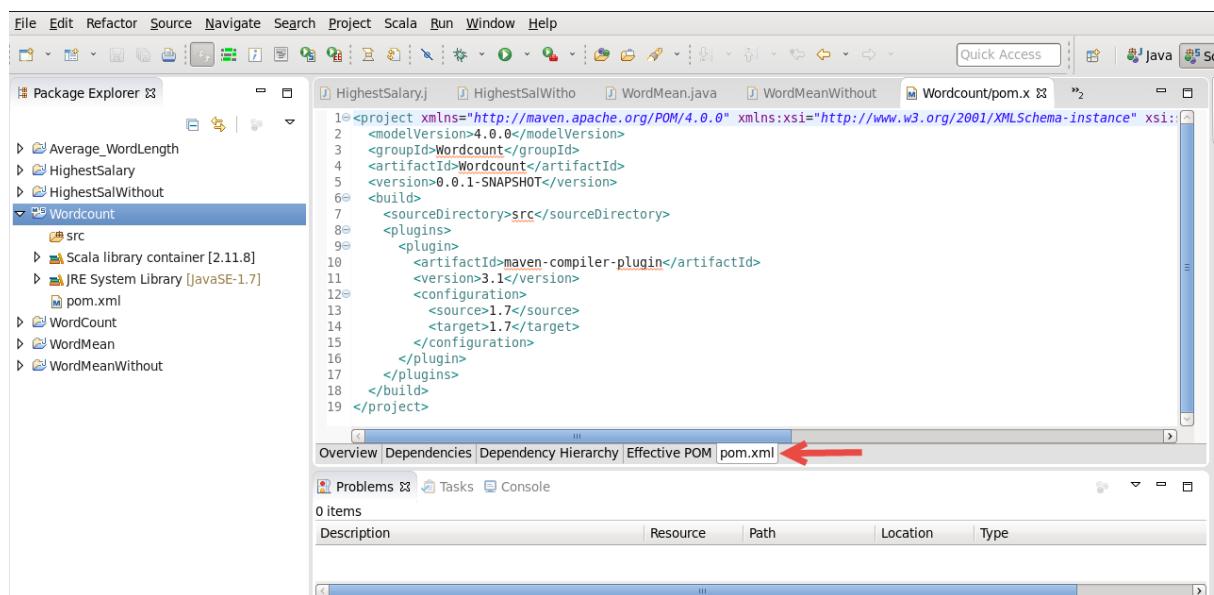
5. Spark 1.6 is supported on Scala 2.10. So set Scala library to version 2.10. *Right Click > Configure > Convert to Maven Project.*

6. Update the maven pom.xml with the following dependences. Keep other content of pom xml as is. Remember to add the below content right under the xml root called <project>.

```

<properties>
  <sparkVersion>1.6.1</sparkVersion>
  <scalaVersion>2.10</scalaVersion>
</properties>
<packaging>jar</packaging>
<dependencies>
  <dependency>
    <groupId>org.apache.spark</groupId>
    <artifactId>spark-core_${scalaVersion}</artifactId>
    <version>${sparkVersion}</version>
    <scope>provided</scope>
  </dependency>
  <dependency>
    <groupId>org.apache.spark</groupId>
    <artifactId>spark-streaming_${scalaVersion}</artifactId>
    <version>${sparkVersion}</version>
  </dependency>
  <dependency>
    <groupId>org.apache.spark</groupId>
    <artifactId>spark-sql_${scalaVersion}</artifactId>
    <version>${sparkVersion}</version>
  </dependency>
  <dependency>
    <groupId>org.apache.spark</groupId>
    <artifactId>spark-mllib_${scalaVersion}</artifactId>
    <version>${sparkVersion}</version>
  </dependency>
  <dependency>
    <groupId>org.apache.spark</groupId>
    <artifactId>spark-graphx_${scalaVersion}</artifactId>
    <version>${sparkVersion}</version>
  </dependency>
</dependencies>

```



7. Now add a new Scala object file "Sample" to project.

```

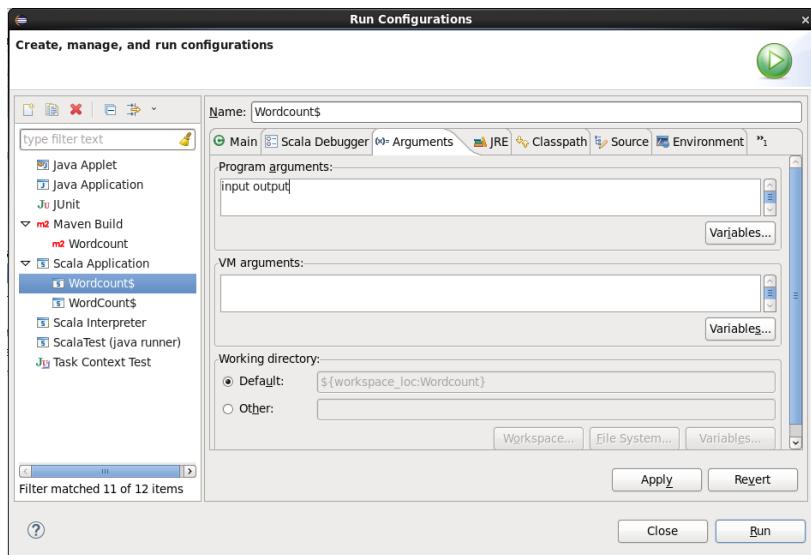
1 import org.apache.spark.{SparkContext, SparkConf}
2
3
4 object WordCount {
5   def main(args: Array[String]): Unit = {
6     if (args.length != 2) {
7       println("Usage: WordCount <input path> <output path>")
8       System.exit(-1)
9     }
10    val conf = new SparkConf()
11
12    conf.setAppName("sample demo")
13    conf.setMaster("local[*]")
14    val sc = new SparkContext(conf)
15    sc.textFile(args(0)) .
16      flatMap(_.split("\\\\W+")) .
17      filter(_.length > 0) .
18      map(_, 1) .
19      reduceByKey(_ + _) .
20      saveAsTextFile(args(1))
21  }
22 }

```

8. Select the file and run it locally from Menu bar Run > Run. This step creates a build configuration. You might fail complaining that input does not exist, that is fine.

9. Right click on the file WordCount.scala > Run As > Run Configurations...  
 > *Select Arguments* tab > In Program Argument, type "input output". This tells that treat input directory under the working directory is the input path and output is the output directory in which the program will save the output.

Now create a input directory and create a test file inside it. Refresh the project. If you find a folder called "output", delete it and delete it before every time you try to run the program again from Eclipse.



10. Now, right click on the project select Run As > Maven Install. This will create a jar file and place under the "target" folder. You have refresh the project to see the jar file.

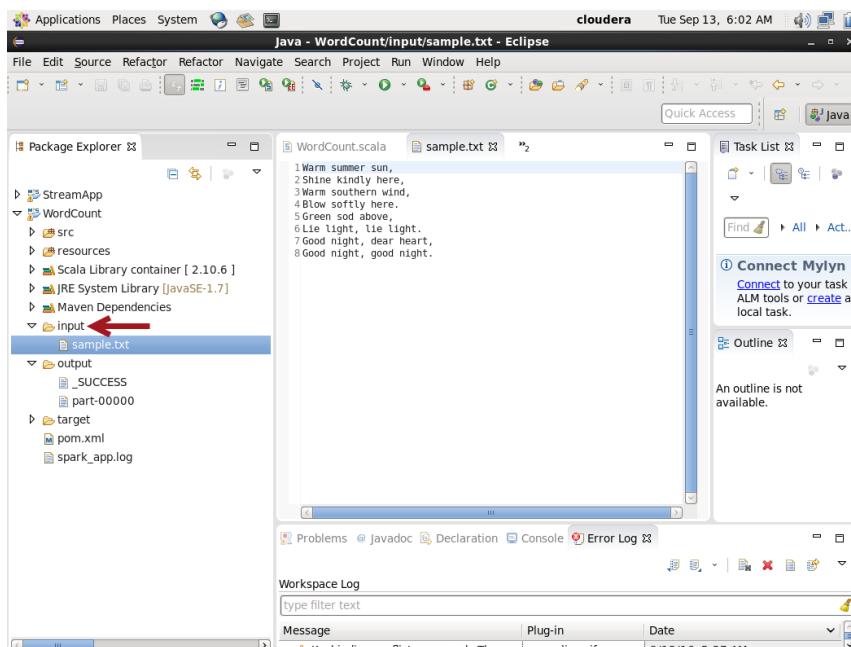
11. Take this jar and submit using spark-submit command

```
cd /home/cloudera/workspace/WordCount/target

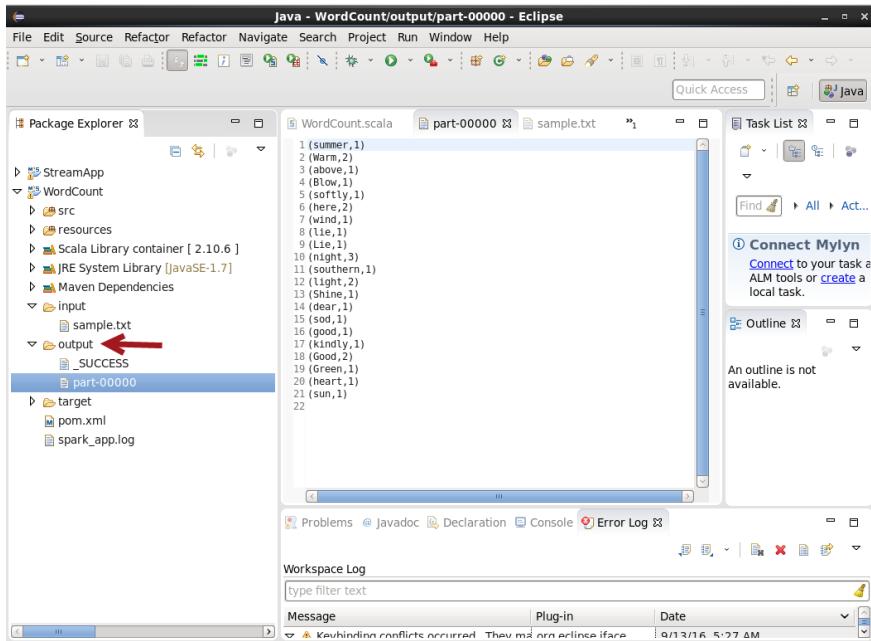
$ spark-submit --class WordCount --master local [*] --jars WordCount-0.0.1-SNAPSHOT.jar
```

```
[cloudera@quickstart target]$ spark-submit --class WordCount --master local[*] WordCount-0.0.1-SNAPSHOT.jar --verbose
Usage: WordCount <input path> <output path>
```

## Input



## Output



Note: Ensure to remove output folder every time you run the program.

### Additional Configuration for Logging

Create directory "resources" and log4j.properties file inside it with the following content.

```
log4j.rootLogger=INFO,console,file

log4j.appender.console=org.apache.log4j.ConsoleAppender
log4j.appender.console.Threshold=WARN
log4j.appender.console.target=System.err
log4j.appender.console.layout=org.apache.log4j.PatternLayout
log4j.appender.console.layout.ConversionPattern=%d{yy/MM/dd HH:mm:ss} %p %c{2}: %m%n

log4j.appender.file=org.apache.log4j.RollingFileAppender
log4j.appender.file.Threshold=INFO
log4j.appender.file.File=spark_app.log
```

```
log4j.appenders.file.MaxFileSize=10MB  
log4j.appenders.file.MaxBackupIndex=10  
log4j.appenders.file.layout=org.apache.log4j.PatternLayout  
log4j.appenders.file.layout.ConversionPattern=%d{yyyy-MM-dd HH:mm:ss} %-  
5p %c{1}:%L - %m%n
```

Right click on the resources folder > Build Path > Use As Source Folder

Now, if you run the above WordCount.scala program again, you will see less debug message. You can control using log handler "file" and "console".

## Additional Exercise

### Analyze Movie Lens Dataset using RDD

**Source:** <http://grouplens.org/datasets/movielens/>

**Objective:** What are the occupations of the users who have rated the iconic movie Titanic as 5?

**Information about Dataset:**

#### **RATINGS FILE DESCRIPTION**

All ratings are contained in the file “**ratings.dat**” and are in the following format:

UserID::MovieID::Rating::Timestamp

UserIDs range between 1 and 6040

MovieIDs range between 1 and 3952

Ratings are made on a 5-star scale (whole-star ratings only)

Timestamp is represented in seconds since the epoch as returned by time(2)

Each user has at least 20 ratings

#### **USERS FILE DESCRIPTION**

User information is in the file “**users.dat**” and is in the following format:

UserID::Gender::Age::Occupation::Zip-code

All demographic information is provided voluntarily by the users and is not checked for accuracy. Only users who have provided some demographic information are included in this data set.

Gender is denoted by a “M” for male and “F” for female

Age is chosen from the following ranges:

1: “Under 18”

18: “18-24”

25: “25-34”

35: “35-44”

45: “45-49”

50: “50-55”

56: "56+"

**Occupation is chosen from the following choices:**

0: "other" or not specified

1: "academic/educator"

2: "artist"

3: "clerical/admin"

4: "college/grad student"

5: "customer service"

6: "doctor/health care"

7: "executive/managerial"

8: "farmer"

9: "homemaker"

10: "K-12 student"

11: "lawyer"

12: "programmer"

13: "retired"

14: "sales/marketing"

15: "scientist"

16: "self-employed"

17: "technician/engineer"

18: "tradesman/craftsman"

19: "unemployed"

20: "writer"

**MOVIES FILE DESCRIPTION**

Movie information is in the file "**movies.dat**" and is in the following format:

MovieID::Title::Genres

Titles are identical to titles provided by the IMDB (including year of release)

Genres are pipe-separated and are selected from the following genres:

Action, Adventure, Animation, Children's, Comedy, Crime, Documentary, Drama, Fantasy, Film-Noir, Horror, Musical, Mystery, Romance, Sci-Fi, Thriller, War, Western

## Steps:

Load the files ‘movies.dat’, ‘users.dat’ and ‘ratings.dat’ into HDFS from local filesystem

1. Create 3 RDD for movies, users and ratings respectively

```
| val movies = sc.textFile("data/movies.dat")
| val users = sc.textFile("data/users.dat")
| val ratings = sc.textFile("data/ratings.dat")
```

2. Print first 10 lines from the movies rdd

```
| movies.take(10).foreach(println)
```

```
16/08/31 23:41:52 INFO scheduler.DAGScheduler: Job 0 finished: take at <console>
:24, took 1.564272 s
1::Toy Story (1995)::Animation|Children's|Comedy
2::Jumanji (1995)::Adventure|Children's|Fantasy
3::Grumpier Old Men (1995)::Comedy|Romance
4::Waiting to Exhale (1995)::Comedy|Drama
5::Father of the Bride Part II (1995)::Comedy
6::Heat (1995)::Action|Crime|Thriller
7::Sabrina (1995)::Comedy|Romance
8::Tom and Huck (1995)::Adventure|Children's
9::Sudden Death (1995)::Action
10::GoldenEye (1995)::Action|Adventure|Thriller
```

3. Find out the id of the movie Titanic

```
| movies.filter(_.contains("Titanic")).collect
```

```
16/08/31 23:42:52 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 1.0, whose t
asks have all completed, from pool
res1: Array[String] = Array(1721::Titanic (1997)::Drama|Romance, 2157::Chamberma
id on the Titanic, The (1998)::Romance, 3403::Raise the Titanic (1980)::Drama|Th
riller, 3404::Titanic (1953)::Action|Drama)
```

4. Print 10 lines from ratings RDD

```
| ratings.take(10).foreach(println)
```

```
16/08/31 23:43:47 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 2.0, whose tasks have all completed, from pool
1::1193::5::978300760
1::661::3::978302109
1::914::3::978301968
1::3408::4::978300275
1::2355::5::978824291
1::1197::3::978302268
1::1287::5::978302039
1::2804::5::978300719
1::594::4::978302268
1::919::4::978301368
```

5. Create a filtered rdd based on the users who have rated Titanic as 1

```
val titanicRating1 = ratings.map(_.split(":")).filter{rec =>
  rec(0).toInt == 1721 && rec(2).toInt == 5}
titanicRating1.take(10).foreach(rec => println(rec.mkString("\t")))
```

```
16/08/31 23:44:48 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 3.0, whose tasks have all completed, from pool
1721    3087      5      974709030
1721    1485      5      974706178
1721    3247      5      974706026
1721    2463      5      974708897
1721    2321      5      974707939
1721    1569      5      974706315
1721    3397      5      974708938
1721    2918      5      974708773
```

6. Convert the previous RDD by assigning user id as key

```
val byUser = titanicRating1.keyBy(_(1))
byUser.take(3)
```

```
res4: Array[(String, Array[String])] = Array((3087,Array(1721, 3087, 5, 974709030)), (1485,Array(1721, 1485, 5, 974706178)), (3247,Array(1721, 3247, 5, 974706026)))
```

7. Show some sample values from users rdd

```
users.take(10).foreach(println)
```

```
16/08/31 23:47:02 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 5.0, whose tasks have all completed, from pool
1::F::1::10::48067
2::M::56::16::70072
3::M::25::15::55117
4::M::45::7::02460
5::M::25::20::55455
6::F::50::9::55117
7::M::35::1::06810
8::M::25::12::11413
9::M::25::17::61614
10::F::35::1::95370
```

#### 8. Create pair RDD pairRddUsers from users rdd using userid as key

```
val pairRddUsers = users.keyBy(_.split("::")(0))
pairRddUsers.take(3)
```

```
16/08/31 23:47:53 INFO scheduler.DAGScheduler: Job 6 finished: take at <console>
:26, took 0.354851 s
res6: Array[(String, String)] = Array((1,1::F::1::10::48067), (2,2::M::56::16::70072), (3,3::M::25::15::55117))
```

#### 9. Join the last Pair RDD to users RDD

```
val rddJoined = byUser.join(pairRddUsers)
rddJoined.take(3)
```

```
16/08/31 23:49:35 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 9.0, whose tasks have all completed, from pool
res7: Array[(String, (Array[String], String))] = Array((3247,(Array(1721, 3247, 5, 974706026),3247::M::35::15::92649)), (3397,(Array(1721, 3397, 5, 974708938),3397::M::25::5::55414)), (1485,(Array(1721, 1485, 5, 974706178),1485::F::25::9::80538)))
```

#### 10. Extract the occupations of the users from the joined RDD

```
rddJoined.map(_.2._2.split("::")(3)).distinct.collect
```

```
16/08/31 23:50:12 INFO scheduler.DAGScheduler: Job 8 finished: collect at <console>:34, took 1.229909 s
res8: Array[String] = Array(7, 15, 5, 9, 12, 13, 1)
```

## Assignment:

**Objective:** Using above transformation and action, Analyse Crime Data from San Francisco Police Department.

### **Steps:**

1. Create an RDD with crime incident data.
2. Take first five records of RDD created and print the same.
3. Reference each field with an index map.
4. Show the first value of an RDD.
5. Check number of partition.
6. Capture the line containing header in RDD.
7. Create an RDD with only records eliminating the header.
8. Find the number of incident records?
9. List all the categories from incident data?
10. Find total number of categories?
11. What are the total number of incidents in each category?
12. Find out on which day each of incidents occurred most?

## Exercise 6: Dataframe and Spark SQL

Two ways to create Spark Dataframes:

1. Create from an existing RDD
2. Create from other data sources

Creating DataFrames from Existing RDD:

1. Infer schema by Reflection
  - a. Convert RDD containing case classes.
  - b. Use when schema is known.
2. Construct schema programmatically
  - a. Schema is dynamic
  - b. Number of fields in case class is more than 22 fields.

### Infer schema by Reflection (case class)

Steps:

1. Import necessary classes
2. Create RDD
3. Define case class
4. Convert RDD into RDD of case objects.
5. Implicitly convert resulting RDD of case objects into DataFrames.
6. Register DataFrame as table.

**Exercise:** Analyze SFPD Crime Dataset using Spark Dataframes.

1. Upload the file ‘crime\_incidents.csv’ from datasets directory on local filesystem to HDFS.
2. Create a base RDD from the file ‘crime\_incidents’.

```
|-----|  
| val rdd = sc.textFile("data/crime_incidents.csv")  
| rdd.take(5).foreach(println)  
|-----|
```

```

IncidntNum,Category,Descript,DayOfWeek,Date,Time,PdDistrict,Resolution,Address,X
,Y,Location
050436712,ASSAULT,BATTERY,Wednesday,04/20/2005 12:00:00 AM,04:00,MISSION,NONE,18
TH ST / CASTRO ST,-122.435002864271,37.7608878061245,"(37.7608878061245, -122.43
5002864271)"
080049078,LARCENY/THEFT,GRAND THEFT FROM A BUILDING,Sunday,01/13/2008 12:00:00 A
M,18:00,PARK,NONE,1100 Block of CLAYTON ST,-122.446837820235,37.7622550270122,"(
37.7622550270122, -122.446837820235)"
130366639,ASSAULT,AGGRAVATED ASSAULT WITH A KNIFE,Sunday,05/05/2013 12:00:00 AM,
04:10,INGLESIDE,"ARREST, BOOKED",0 Block of SGTJOHNVYOUNG LN,-122.444707063455,3
7.7249307267936,"(37.7249307267936, -122.444707063455)"
030810835,DRIVING UNDER THE INFLUENCE,DRIVING WHILE UNDER THE INFLUENCE OF ALCOH
OL,Tuesday,07/08/2003 12:00:00 AM,01:00,SOUTHERN,"ARREST, BOOKED",MASON ST / TUR
K ST,-122.408953598286,37.7832878735491,"(37.7832878735491, -122.408953598286)"

```

### 3. Cache the RDD since we are going to process it several times

```

-----[rdd.cache
rdd.count]-----
```

```

16/09/01 02:32:12 INFO scheduler.DAGScheduler: Job 1 finished: count at <console
>:24, took 13.574773 s
res2: Long = 1888568

```

### 4. Create a filtered RDD by removing the headers

```

-----[val data = rdd.filter(line => !line.startsWith("IncidntNum"))
data.take(5).foreach(println)]-----
```

```

16/09/01 02:33:04 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 2.0, whose t
asks have all completed, from pool
050436712,ASSAULT,BATTERY,Wednesday,04/20/2005 12:00:00 AM,04:00,MISSION,NONE,18
TH ST / CASTRO ST,-122.435002864271,37.7608878061245,"(37.7608878061245, -122.43
5002864271)"
080049078,LARCENY/THEFT,GRAND THEFT FROM A BUILDING,Sunday,01/13/2008 12:00:00 A
M,18:00,PARK,NONE,1100 Block of CLAYTON ST,-122.446837820235,37.7622550270122,"(
37.7622550270122, -122.446837820235)"
130366639,ASSAULT,AGGRAVATED ASSAULT WITH A KNIFE,Sunday,05/05/2013 12:00:00 AM,
04:10,INGLESIDE,"ARREST, BOOKED",0 Block of SGTJOHNVYOUNG LN,-122.444707063455,3
7.7249307267936,"(37.7249307267936, -122.444707063455)"
030810835,DRIVING UNDER THE INFLUENCE,DRIVING WHILE UNDER THE INFLUENCE OF ALCOH
OL,Tuesday,07/08/2003 12:00:00 AM,01:00,SOUTHERN,"ARREST, BOOKED",MASON ST / TUR
K ST,-122.408953598286,37.7832878735491,"(37.7832878735491, -122.408953598286)"
130839567,OTHER OFFENSES,TRAFFIC VIOLATION ARREST,Friday,10/04/2013 12:00:00 AM,
20:53,TENDERLOIN,"ARREST, BOOKED",TURK ST / LEAVENWORTH ST,-122.414056291891,37.
7827931071006,"(37.7827931071006, -122.414056291891)"
```

5. Parse the data using regex. Regex allows to handle fields that contain the field separators comma

```

val pattern = """([^\"]*)"" | (?<=,|^)([^,]*)(?=,|$)""".r
val dataParsed = data.map(pattern.findAllIn(_).toArray)
dataParsed.take(5).foreach(rec => println(rec.mkString(", ")))

```

```

16/09/01 02:34:42 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 3.0, whose tasks have all completed, from pool
050436712,ASSAULT,BATTERY,Wednesday,04/20/2005 12:00:00 AM,04:00,MISSION,NONE,18TH ST / CASTRO ST,-122.435002864271,37.7608878061245,"(37.7608878061245, -122.435002864271)"
080049078,LARCENY/THEFT,GRAND THEFT FROM A BUILDING,Sunday,01/13/2008 12:00:00 AM,18:00,PARK,NONE,1100 Block of CLAYTON ST,-122.446837820235,37.7622550270122,"(37.7622550270122, -122.446837820235)"
130366639,ASSAULT,AGGRAVATED ASSAULT WITH A KNIFE,Sunday,05/05/2013 12:00:00 AM,04:10,INGLESIDE,"ARREST, BOOKED",0 Block of SGTJOHNVYOUNG LN,-122.444707063455,37.7249307267936,"(37.7249307267936, -122.444707063455)"
030810835,DRIVING UNDER THE INFLUENCE,DRIVING WHILE UNDER THE INFLUENCE OF ALCOHOL,Tuesday,07/08/2003 12:00:00 AM,01:00,SOUTHERN,"ARREST, BOOKED",MASON ST / TURK ST,-122.408953598286,37.7832878735491,"(37.7832878735491, -122.408953598286)"
130839567,OTHER OFFENSES,TRAFFIC VIOLATION ARREST,Friday,10/04/2013 12:00:00 AM,20:53,TENDERLOIN,"ARREST, BOOKED",TURK ST / LEAVENWORTH ST,-122.414056291891,37.7827931071006,"(37.7827931071006, -122.414056291891)"

```

### Create Dataframe using case class objects

6. Create a case class that we will encapsulate the incident records into

```

case class Incident(
    IncidntNum:String,
    Category:String,
    Description:String,
    DayOfWeek:String,
    DateOfIncident:String,
    TimeOfIncident:String,
    PdDistrict:String,
    Resolution:String,
    Address:String,
    X:String,
    Y:String,
    Location:String)

```

## 7. Create a new rdd from dataParsed by converting the records into case class objects

```
val incidentCaseObjects = dataParsed.map(rec => Incident(rec(0),  
rec(1), rec(2), rec(3), rec(4), rec(5), rec(6), rec(7), rec(8),  
rec(9), rec(10), rec(11)))  
  
incidentCaseObjects.take(10).foreach(println)
```

## 8. Convert the rdd of the case class objects into a dataframe

```
val incidentsDF = incidentCaseObjects.toDF  
  
incidentsDF.show
```

IncidntNum	Category	Description	DayOfWeek	DateOfIncident	TimeOfIncident	PdDistrict	Resolution	Address	X	Y	Location
050436712	ASSAULT	BATTERY	Wednesday	04/28/2008	12:08:...	04:00	MISSION	NONE 18TH ST / CASTRO ST -122,435082864271 37.7608878061245 "37.7608878061245			
088649078	LARCENY/THEFT	GRAND THEFT FROM ...	Sunday	01/13/2008	12:00:...	18:00	PARK	NONE 1100 Block of CLA... -122,446837820235 37.7622550270122 "37.7622550270122			
130366639	ASSAULT	AGGRAVATED ASSAULT...	Sunday	05/05/2013	12:08:...	04:10	INGLESIDE	"ARREST, BOOKED"0 Block of SGTJMH... -122,444707063455 37.7249307267936 "37.7249307267936			
030810835	DRIVING UNDER THE...	DRIVING WHILE UND...	Tuesday	07/08/2003	12:08:...	01:00	SOUTHERN	"ARREST, BOOKED"1 MASON ST / TURK ST -122,408953598296 37.7832878735491 "37.7832878735491			
130839567	OTHER OFFENSES	TRAFFIC VIOLATION...	Friday	10/04/2013	12:08:...	20:53	TEHOMERLOIN	"ARREST, BOOKED"1 TURK ST / LEAVENN... -122,414065291891 37.7827931071066 "37.7827931071066			
078838580	BURGLARY	BURGLARY OF APAR...	Tuesday	08/14/2008	12:08:...	07:00	NORTHERN	NONE 3100 Block of FRA... -122,426730544229 37.803467499672 "37.803467499672			
080233182	DRUG/NARCOTIC	POSSESSION OF MAR...	Tuesday	03/04/2008	12:08:...	14:23	INGLESIDE	"ARREST, CITED"MISSION ST / PERS... -122,43597721703 37.7231288306727 "37.7231288306727			
060711805	OTHER OFFENSES	DRIVER'S LICENSE, ...	Wednesday	07/05/2006	12:08:...	15:50	INGLESIDE	"ARREST, CITED"2300 Block of SAN... -122,447241159611 37.720577971255 "37.720577971255			
040062593	LARCENY/THEFT	GRAND THEFT FROM ...	Wednesday	12/10/2003	12:08:...	09:30	INGLESIDE	NONE 0 Block of MOFFIT... -122,432787775164 37.7371566745272 "37.7371566745272			
110065182	NON-CRIMINAL	STAY AWAY OR COU...	Monday	01/17/2011	12:08:...	15:35	INGLESIDE	NONE 600 Block of CAMP... -122,40876107232 37.715900951041 "37.715900951041			
0606027038	LARCENY/THEFT	GRAND THEFT FROM ...	Saturday	01/07/2006	12:08:...	22:00	NORTHERN	NONE 100 Block of HEMI... -122,428815262607 37.7872368476925 "37.7872368476925			
110929398	LARCENY/THEFT	PETTY THEFT BICYCLE	Sunday	11/13/2011	12:08:...	18:00	MISSION	NONE 3000 Block of 22N... -122,416014728293 37.7555464259299 "37.7555464259299			

## 9. Use Dataframe API to run analytics on dataframe

```
incidentsDF.groupBy($"Category").count().sort($"count".desc).show
```

```
+-----+-----+
|      Category| count|
+-----+-----+
| LARCENY/THEFT|385774|
| OTHER OFFENSES|269250|
| NON-CRIMINAL|200942|
| ASSAULT|165324|
| VEHICLE THEFT|114258|
| DRUG/NARCOTIC|111436|
| VANDALISM| 96350|
| WARRANTS| 89782|
| BURGLARY| 78968|
| SUSPICIOUS OCC| 67788|
| MISSING PERSON| 55584|
| ROBBERY| 48713|
| FRAUD| 36004|
| FORGERY/COUNTERFE...| 21804|
```

10. Cache the dataframe and compare the size of the dataframe in memory with that of the rdd cached earlier

```
incidentsDF.cache  
incidentsDF.count
```

```
res91: incidentsDF.type = [IncidentNum: string, Category: string, Description: string, DayOfWeek: string, DateOfIncide  
res92: Long = 1888567
```

11. Register the data frame as temporary table in Spark SQL

```
incidentsDF.registerTempTable("incidents")
```

12. Run SQL query on top of newly created table

```
sqlContext.sql("select * from incidents").show
```

[IncidentNum]	Category	Description	DayOfWeek	DateOfIncident	TimeOfIncident	PdDistrict	Resolution	Address	X	Y	Location	
056436712	ASSAULT	BATTERY	Wednesday	04/20/2005	12:00:00	04:00	MISSION	NONE	18TH ST / CASTRO ST	-122.439002864271	37.7686878861245	(37.7686878861245, -122.439002864271)
060043078	LARCENY/THEFT	GRAND THEFT FROM A BUILDING	Sunday	01/13/2008	12:00:00	18:00	PARK	NONE	1100 Block of CLAYTON ST	-122.44683782935	37.762255027012	(37.762255027012, -122.44683782935)
130366391	ASSAULT	AGGRAVATED ASSAULT WITH A KNIFE	Sunday	05/05/2013	12:00:00	04:10	INGLESIDE	"ARREST, BOOKED"	10 Block of SGT JONH YOUNG LN	-122.444707065455	37.724939736793	(37.724939736793, -122.444707065455)
030818835	DRIVING UNDER THE INFLUENCE	DRIVING WHILE UNLAWFUL	Tuesday	07/09/2003	12:00:00	01:00	SOUTHERN	"ARREST, BOOKED"	MASON ST / TURK ST	-122.408953598286	37.783287875491	(37.783287875491, -122.408953598286)
130809567	OTHER OFFENSES	TRAFFIC VIOLATION	Friday	10/04/2007	12:00:00	20:53	TENDERLOIN	"ARREST, BOOKED"	TURK ST / LEAVENWORTH ST	-122.414056291891	37.782793107100	(37.782793107100, -122.414056291891)
076838580	BURGLARY	BURGLARY OF APARTMENT HOUSE, UNLAWFUL ENTRY	Tuesday	08/14/2007	12:00:00	07:00	NORTHERN	NONE	3100 Block of FRANKLIN ST	-122.426730544229	37.8034674969672	(37.8034674969672, -122.426730544229)
088223102	DRUG/NARCOTIC	POSSESSION OF MARIJUANA	Tuesday	03/04/2008	12:00:00	14:23	INGLESIDE	"ARREST, CITED"	MISSION ST / PERSIA AV	-122.435977217671	37.7231288366727	(37.7231288366727, -122.435977217671)
110711805	OTHER OFFENSES	DRIVERS LICENSE, SUSPENDED OR REVOKED	Wednesday	07/05/2006	12:00:00	15:50	INGLESIDE	"ARREST, CITED"	2300 Block of SAN JOSE AV	-122.4211596131	37.720157797125	(37.720157797125, -122.4211596131)
048062593	LARCENY/THEFT	GRAND THEFT FROM A BUILDING	Wednesday	12/10/2003	12:00:00	09:30	INGLESIDE	NONE	0 Block of MOFFIT ST	-122.432787775164	37.7371566745272	(37.7371566745272, -122.432787775164)
110851822	NON-CRIMINAL	"STAY AWAY OR CO." ORDER	Monday	01/17/2011	12:00:00	15:35	INGLESIDE	NONE	600 Block of CAMP ST	-122.408761072232	37.7159008951041	(37.7159008951041, -122.408761072232)
066027038	LARCENY/THEFT	GRAND THEFT FROM A BUILDING	Saturday	01/07/2006	12:00:00	22:00	NORTHERN	NONE	100 Block of HEMINGWAY ST	-122.426815202607	37.7872360476925	(37.7872360476925, -122.426815202607)

13. And you can also run SQL statement through SQL interpreter in Zeppelin

```
%sql select * from incidents limit 10
```

IncidentNum	Category	Description	DayOfWeek	DateOfIncident	TimeOfIncident	PdDistrict	Resolution	Address	X	Y	Location	
50,436,712	ASSAULT	BATTERY	Wednesday	04/20/2005	12:00:00 AM	04:00	MISSION	NONE	18TH ST / CASTRO ST	-122.435	37.76089	(37.7608878061245, -122.435)
80,049,078	LARCENY/THEFT	GRAND THEFT FROM A BUILDING	Sunday	01/13/2008	12:00:00 AM	18:00	PARK	NONE	1100 Block of CLAYTON ST	-122.44683	37.76226	(37.762255027012, -122.44683782935)
130,366,639	ASSAULT	AGGRAVATED ASSAULT WITH A KNIFE	Sunday	05/05/2013	12:00:00 AM	04:10	INGLESIDE	"ARREST, BOOKED"	0 Block of SGT JONH YOUNG LN	-122.44471	37.72493	(37.724939736793, -122.444707065455)
30,810,835	DRIVING UNDER THE INFLUENCE	DRIVING WHILE UNDER THE INFLUENCE OF ALCOHOL	Tuesday	07/09/2003	12:00:00 AM	01:00	SOUTHERN	"ARREST, BOOKED"	MASON ST / TURK ST	-122.40895	37.78329	(37.783287875491, -122.408953598286)
130,839,567	OTHER OFFENSES	TRAFFIC VIOLATION ARREST	Friday	10/04/2013	12:00:00 AM	20:53	TENDERLOIN	"ARREST, BOOKED"	TURK ST / LEAVENWORTH ST	-122.41406	37.78279	(37.782793107100, -122.414056291891)
70,838,580	BURGLARY	"BURGLARY OF APARTMENT HOUSE, UNLAWFUL ENTRY"	Tuesday	08/14/2007	12:00:00 AM	07:00	NORTHERN	NONE	3100 Block of FRANKLIN ST	-122.42673	37.80347	(37.8034674969672, -122.426730544229)
80,233,102	DRUG/NARCOTIC	POSSESSION OF MARIJUANA	Tuesday	03/04/2008	12:00:00 AM	14:23	INGLESIDE	"ARREST, CITED"	MISSION ST / PERSIA AV	-122.43598	37.72313	(37.7231288366727, -122.435987775164)
60,711,805	OTHER OFFENSES	DRIVERS LICENSE, SUSPENDED OR REVOKED	Wednesday	07/05/2006	12:00:00 AM	15:50	INGLESIDE	"ARREST, CITED"	2300 Block of SAN JOSE AV	-122.44724	37.72016	(37.720157797125, -122.447247775164)

```
%sql select category, count(*) `count` from incidents group by category order by `count` desc limit 10
```

category	count
LARCENY/THEFT	385,774
OTHER OFFENSES	269,250
NON-CRIMINAL	200,942
ASSAULT	165,324
VEHICLE THEFT	114,258
DRUG/NARCOTIC	111,436
VANDALISM	96,350
WARRANTS	89,782
BURGLARY	78,968

## Programmatically

### Steps:

1. Create an RDD of Rows from the original RDD
2. Create a schema represented by a StructType matching the structure of Rows in the RDD created in step1.
3. Create DataFrame by applying schema to Row RDD.

## Create Dataframe Dynamically (...without case class objects)

### 1. Create a schema programmatically

```
import org.apache.spark.sql.types.{StructField, StructType,
StringType}

val schema = StructType(Array(
    StructField("IncidentNum", StringType, false),
    StructField("Category", StringType, false),
    StructField("Description", StringType, false),
    StructField("DayOfWeek", StringType, false),
    StructField("DateOfIncident", StringType, false),
    StructField("TimeOfIncident", StringType, false),
    StructField("PdDistrict", StringType, false),
    StructField("Resolution", StringType, false),
    StructField("Address", StringType, false),
    StructField("X", StringType, false),
    StructField("Y", StringType, false),
    StructField("Location", StringType, false)
))
```

### 2. Convert the data to RDD of Row objects

```
import org.apache.spark.sql.Row

val incidentRowObjects = dataParsed.map(rec => Row(rec(0), rec(1),
rec(2), rec(3), rec(4), rec(5), rec(6), rec(7), rec(8), rec(9),
rec(10), rec(11)))
```

### 3. Create the dataframe from RDD of Row objects and dynamic schema

```
val incidentsDF = sqlContext.createDataFrame(incidentRowObjects,
schema)
```

- Use dataframe api and further, register the data frame as Spark SQL table and run SQL statements on top of it

```
incidentsDF.show
```

[IncidentNum]	Category	Description DayOfWeek	DateOfIncident TimeOfIncident	PdDistrict	Resolution	Address	X	Y	Location
056436712	ASSAULT	BATTERY Wednesday 04/20/2005 12:00:...	04:00	MISSION	NONE	18TH ST / CASTRO ST	-122.435062864271	37.7608878061245	(37.7608878061245, -122.435062864271)
088049078	LARCENY/THEFT	GRAND THEFT FROM ... Sunday 01/13/2008 12:00:...	18:00	PARK	NONE	1100 Block of CLA...	-122.446837826235	37.7622550276122	(37.7622550276122, -122.446837826235)
138366639	ASSAULT	AGGRAVATED ASSAUL... Sunday 05/05/2013 12:00:...	04:10	INGLESIDE	"ARREST, BOOKED"	0 Block of SGTJOH...	-122.444767063455	37.7249307267936	(37.7249307267936, -122.444767063455)
036810835	DRIVING UNDER THE...	DRIVING WHILE UND... Tuesday 07/08/2003 12:00:...	01:00	SOUTHERN	"ARREST, BOOKED"	MASON ST / TURK ST	-122.408953598286	37.7832878735491	(37.7832878735491, -122.408953598286)
138639567	OTHER OFFENSES	TRAFFIC VIOLATION... Friday 10/04/2013 12:00:...	20:53	TENDERLOIN	"ARREST, BOOKED"	TURK ST / LEAVENW...	-122.414656219891	37.7827931071006	(37.7827931071006, -122.414656219891)
076838580	BURGLARY	BURGLARY OF APAR... Tuesday 08/14/2007 12:00:...	07:00	NORTHERN	NONE	3100 Block of FRA...	-122.426730544229	37.8034674969672	(37.8034674969672, -122.426730544229)
086233102	DRUG/NARCOTIC	POSSESSION OF MAR... Tuesday 03/04/2008 12:00:...	14:23	INGLESIDE	"ARREST, CITED"	MISSION ST / FERS...	-122.43597721703	37.7231288306722	(37.7231288306722, -122.43597721703)
066711805	OTHER OFFENSES	DRIVERS LICENSE,... Wednesday 07/05/2006 12:00:...	15:50	INGLESIDE	"ARREST, CITED"	2300 Block of SAN...	-122.447241159611	37.7201577971255	(37.7201577971255, -122.447241159611)
040662593	LARCENY/THEFT	GRAND THEFT FROM ... Wednesday 12/10/2003 12:00:...	09:30	INGLESIDE	NONE	0 Block of MOFFIT...	-122.432787775164	37.7371566745272	(37.7371566745272, -122.432787775164)
110651822	NON-CRIMINAL	STAY AWAY OR COU... Monday 01/17/2011 12:00:...	15:35	INGLESIDE	NONE	600 Block of CAMP...	-122.4089761072232	37.7159080951041	(37.7159080951041, -122.4089761072232)
066027038	LARCENY/THEFT	GRAND THEFT FROM ... Saturday 01/07/2006 12:00:...	22:00	NORTHERN	NONE	100 Block of HEML...	-122.420815262697	37.7872360476925	(37.7872360476925, -122.420815262697)

## UDF in Spark Dataframe

There are two type:

- UDF for Dataframe API
- UDF for SQL expressions

### UDF for Dataframe API

- Parse string DateOfIncident field into a time stamp

```
import java.text.SimpleDateFormat
import java.util.Locale
import java.sql.Timestamp // Spark SQL does not support java.util.Date

def sf = new SimpleDateFormat( "MM/dd/yyyy HH:mm:ss aaa",
Locale.ENGLISH) //Sample value: 10/04/2013 12:00:00 AM
sf.setLenient(true)

def parseDatetimeString(dateStr: String): Timestamp = {
    new Timestamp(sf.parse(dateStr).getTime)
}
```

## 2. Register the UDF for Dataframe API

```
import org.apache.spark.sql.functions.udf  
def toTimestamp = udf(parseDatetimeString _)
```

## 3. Use the UDF you just created

```
incidentsDF.select($"DateOfIncident",  
toTimestamp($"DateOfIncident")).show(10, false)
```

```
+-----+  
|DateOfIncident |UDF(DateOfIncident) |  
+-----+  
|04/20/2005 12:00:00 AM|2005-04-20 12:00:00.0|  
|01/13/2008 12:00:00 AM|2008-01-13 12:00:00.0|  
|05/05/2013 12:00:00 AM|2013-05-05 12:00:00.0|  
|07/08/2003 12:00:00 AM|2003-07-08 12:00:00.0|  
|10/04/2013 12:00:00 AM|2013-10-04 12:00:00.0|  
|08/14/2007 12:00:00 AM|2007-08-14 12:00:00.0|  
|03/04/2008 12:00:00 AM|2008-03-04 12:00:00.0|  
|07/05/2006 12:00:00 AM|2006-07-05 12:00:00.0|  
|12/10/2003 12:00:00 AM|2003-12-10 12:00:00.0|  
|01/17/2011 12:00:00 AM|2011-01-17 12:00:00.0|  
+-----+  
only showing top 10 rows
```

## UDF for SQL Expression

### 1. Register UDF for SQL expression

```
sqlContext.udf.register("toTimestamp", parseDatetimeString _)
```

### 2. Use UDF in SQL statements that you just created

```
sqlContext.sql("select t.IncidentNum, t.DateOfIncident,  
toTimestamp(DateOfIncident) `Timestamp` from incidents t").show(10,  
false)
```

### 3. Use them in sql interpreter as well using zeppelin

```
%sql select t.IncidntNum, t.DateOfIncident,  
toTimestamp(DateOfIncident) `Timestamp` from incidents t limit 10
```

IncidntNum	DateOfIncident	Timestamp
50,436,712	04/20/2005 12:00:00 AM	2005-04-20 12:00:00.0
80,049,078	01/13/2008 12:00:00 AM	2008-01-13 12:00:00.0
130,366,639	05/05/2013 12:00:00 AM	2013-05-05 12:00:00.0
30,810,835	07/08/2003 12:00:00 AM	2003-07-08 12:00:00.0
130,839,567	10/04/2013 12:00:00 AM	2013-10-04 12:00:00.0
70,838,580	08/14/2007 12:00:00 AM	2007-08-14 12:00:00.0
80,233,102	03/04/2008 12:00:00 AM	2008-03-04 12:00:00.0
60,711,805	07/05/2006 12:00:00 AM	2006-07-05 12:00:00.0

### Save the Dataframe to filesystem

1. Default format is ‘parquet’

```
incidentsDF.write.save("/data/incidents")
```

2. To save in ‘JSON’ format

```
incidentsDF.write.format("json").save("/data/incidents_json")
```

### Creating Dataframe from certain file format

```
val incidents = sqlContext.read.format("com.databricks.spark.csv").  
options(Map("header" -> "true", "inferSchema" -> "true")).  
load("/data/Map__Crime_Incidents_-_from_1_Jan_2003.csv")  
incidents.show
```

IncidntNum	DateOfIncident	Timestamp
50,436,712	04/20/2005 12:00:00 AM	2005-04-20 12:00:00.0
80,049,078	01/13/2008 12:00:00 AM	2008-01-13 12:00:00.0
130,366,639	05/05/2013 12:00:00 AM	2013-05-05 12:00:00.0
30,810,835	07/08/2003 12:00:00 AM	2003-07-08 12:00:00.0
130,839,567	10/04/2013 12:00:00 AM	2013-10-04 12:00:00.0
70,838,580	08/14/2007 12:00:00 AM	2007-08-14 12:00:00.0
80,233,102	03/04/2008 12:00:00 AM	2008-03-04 12:00:00.0
60,711,805	07/05/2006 12:00:00 AM	2006-07-05 12:00:00.0

## Exercise 7: Spark Streaming

Let's start with simple streaming example, to count the number of words in text data received from a data server listening on a TCP socket.

Create a DStream that represents streaming data from a TCP source, specified as hostname (e.g. localhost) and port (e.g. 9999).

Create a Scala project with below inputs:

```
-----  
import org.apache.spark.SparkConf  
  
import org.apache.spark.SparkContext._  
  
import org.apache.spark.streaming.StreamingContext  
  
import org.apache.spark.streaming.Seconds  
  
import org.apache.spark.storage.StorageLevel;  
  
import org.apache.spark.streaming.StreamingContext._  
  
object SparkStreamReceiver {  
  
    def main(args: Array[String]): Unit = {  
  
        var port:Int = 9999;  
  
        val conf = new SparkConf()  
  
            .setMaster("local[*]")  
  
            .setAppName("Simple spark streaming")  
  
  
        val ssc = new StreamingContext(conf, Seconds(3));  
  
        val dstream = ssc.socketTextStream("localhost", port,  
StorageLevel.MEMORY_ONLY);  
  
        val words = dstream.flatMap(_.split(" "))  
        val pairs = words.map(word => (word, 1))  
        val wordcounts = pairs.reduceByKey(_ + _)  
        wordcounts.print();  
        ssc.start();  
        ssc.awaitTermination();  
    }  
}  
-----
```

Export the jar file i.e. StreamApp.jar

```
[cloudera@quickstart workspace]$ ls -ltr
total 28
drwxrwxr-x 5 cloudera cloudera 4096 Sep 13 06:23 Wordcount
-rw-rw-r-- 1 cloudera cloudera 6055 Sep 13 06:40 wordcount.jar
drwxrwxr-x 7 cloudera cloudera 4096 Sep 13 07:33 StreamApp
drwxrwxr-x 8 cloudera cloudera 4096 Sep 13 07:33 WordCount
-rw-rw-r-- 1 cloudera cloudera 6920 Sep 13 07:45 StreamApp.jar
[cloudera@quickstart workspace]$ █
```

Run the program using spark-submit

Run Netcat command and open another terminal to submit the job

```
| nc -nlk 9999
| -----
```

```
[cloudera@quickstart workspace]$ nc -nlk 9999
hello world
hello hadoop
█
```

```
| spark-submit --class SparkStreamReceiver --master local[3] --verbose
| StreamApp.jar
| -----
```

Any lines typed in the terminal running the netcat server will be counted and printed on screen every second. It will look something like the following

```
16/09/13 08:19:03 INFO executor.Executor: Finished task 0.0 in stage 4.0 (TID 3)
. 1312 bytes result sent to driver
16/09/13 08:19:03 INFO scheduler.TaskSetManager: Finished task 0.0 in stage 4.0
(TID 3) in 10 ms on localhost (1/1)
16/09/13 08:19:03 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 4.0, whose t
asks have all completed, from pool
16/09/13 08:19:03 INFO scheduler.DAGScheduler: ResultStage 4 (print at SparkStre
amReceiver.scala:24) finished in 0.012 s
16/09/13 08:19:03 INFO scheduler.DAGScheduler: Job 2 finished: print at SparkStr
eamReceiver.scala:24, took 0.033085 s
-----
Time: 1473779943000 ms
-----
(hello,2)
(world,1)
(hadoop,1)
```

## Exercise 8: SparkML

### Predicting power grid demand using Spark ML

**Data Source:** <https://archive.ics.uci.edu/ml/datasets/Combined+Cycle+Power+Plant>

Features consist of hourly average ambient variables

Temperature (T) in the range 1.81°C and 37.11°C,

Ambient Pressure (AP) in the range 992.89-1033.30 milibar,

Relative Humidity (RH) in the range 25.56% to 100.16%

Exhaust Vacuum (V) in the range 25.36-81.56 cm Hg

Net hourly electrical energy output (EP) 420.26-495.76 MW

The averages are taken from various sensors located around the plant that record the ambient variables every second. The variables are given without normalization.

#### Dataset Information:

The dataset contains 9568 data points collected from a Combined Cycle Power Plant over 6 years (2006-2011), when the power plant was set to work with full load. Features consist of hourly average ambient variables Temperature (T), Ambient Pressure (AP), Relative Humidity (RH) and Exhaust Vacuum (V) to predict the net hourly electrical energy output (EP) of the plant.

A combined cycle power plant (CCPP) is composed of gas turbines (GT), steam turbines (ST) and heat recovery steam generators. In a CCPP, the electricity is generated by gas and steam turbines, which are combined in one cycle, and is transferred from one turbine to another. While the Vacuum is collected from and has effect on the Steam Turbine, the other three of the ambient variables effect the GT performance.

#### Steps:

1. Load the data

```
val df = sqlContext.read.format("csv")
    .option("header", "true")
    .option("inferSchema", "true")
    .load("/data/Combined_Cycle_Power_Plant/")
```

```
| df.printSchema
```

```
| df: org.apache.spark.sql.DataFrame = [AT: double, V: double, AP: double, RH: double, PE: double]
| root
|   |-- AT: double (nullable = true)
|   |-- V: double (nullable = true)
|   |-- AP: double (nullable = true)
|   |-- RH: double (nullable = true)
|   |-- PE: double (nullable = true)
```

2. Cache the dataframe, as it will be used many time.

```
| df.cache
```

```
| df.count
```

```
| df.show(5)
```

```
| res59: df.type = [AT: double, V: double, AP: double, RH: double, PE: double]
| res60: Long = 47840
+---+---+---+---+
|   AT|    V|    AP|    RH|    PE|
+---+---+---+---+
| 14.96| 41.76| 1024.07| 73.17| 463.26|
| 25.18| 62.96| 1020.04| 59.08| 444.37|
|  5.11|  39.4| 1012.16| 92.14| 488.56|
| 20.86| 57.32| 1010.24| 76.64| 446.48|
| 10.82|  37.5| 1009.23| 96.62| 473.9|
+---+---+---+---+
only showing top 5 rows
```

3. Register the dataframe as table

```
| df.registerTempTable("pp")
```

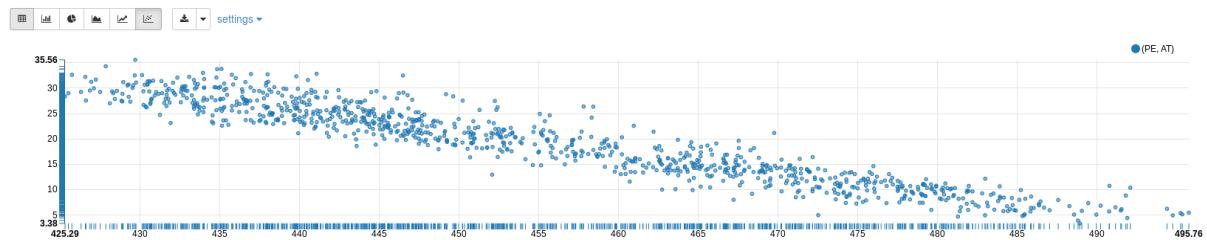
4. Query the table using SQL interpreter from zeeplein.

```
| %sql
```

```
| SELECT PE,AT from pp
```

PE	AT
463.26	14.96
444.37	25.18
488.56	5.11
446.48	20.86
473.9	10.82
443.67	26.27
467.35	15.89
478.42	9.48

To view the output graphically



## 5. Query the table for analysis between PE vs V

```
%sql
SELECT PE, V FROM pp
```

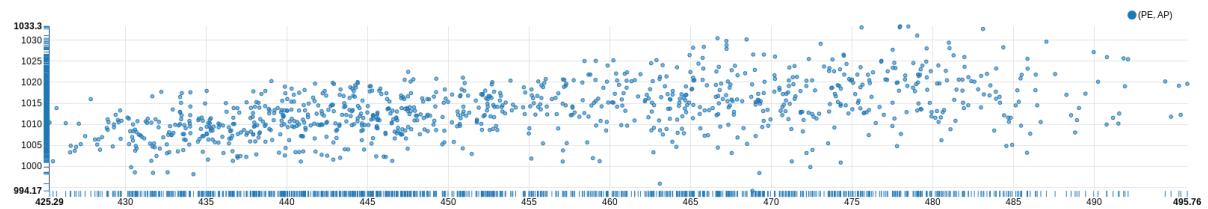
PE	V
463.26	41.76
444.37	62.96
488.56	39.4
446.48	57.32
473.9	37.5
443.67	59.44
467.35	43.96
478.42	44.71
475.98	45

Results are limited by 1000.



## 6. Analysis between PE and AP

```
%sql
SELECT PE, AP FROM pp
```



## 7. Analysis between PE and RH

```
%sql
SELECT PE, RH FROM pp
```

PE	RH
463.26	73.17
444.37	59.08
488.56	92.14
446.48	76.64
473.9	96.62
443.67	58.77
467.35	75.24

## 8. Data preparation

- VectorAssembler is a transformer that combines a given list of columns into a single vector column.
- Convert the `power\_plant` SQL table into a DataFrame named `dataset`

- Set the vectorizer's input columns to a list of the four columns of the input DataFrame: `["AT", "V", "AP", "RH"]`
- Set the vectorizer's output column name to `"features"`

```

import org.apache.spark.ml.feature.VectorAssembler
val vectorizer = new VectorAssembler()
vectorizer.setInputCols(Array("AT", "V", "AP", "RH"))
vectorizer.setOutputCol("features")

```

```

vectorizer: org.apache.spark.ml.feature.VectorAssembler = vecAssembler_c81487e45a35
res66: vectorizer.type = vecAssembler_c81487e45a35
res67: vectorizer.type = vecAssembler_c81487e45a35

```

## 9. Data Modeling

- Split the data into training and test datasets, i.e. 20% for testing and 80% for training.

### Objective:

We need a way of evaluating how well our linear regression model predicts power output as a function of input parameters. We can do this by splitting up our initial data set into a \_Training Set\_ used to train our model and a \_Test Set\_ used to evaluate the model's performance in giving predictions.

```

var Array(testingSet, trainingSet) = df.randomSplit(Array(0.20,
0.80), 0L)
testingSet.count
trainingSet.count

```

```

testingSet: org.apache.spark.sql.DataFrame = [AT: double, V: double, AP: double, RH: double, PE: double]
trainingSet: org.apache.spark.sql.DataFrame = [AT: double, V: double, AP: double, RH: double, PE: double]
res69: Long = 9547
res70: Long = 38293

```

## Linear Regression Model

We'll create a Linear Regression Model and use the built in help to identify how to train it.

- Initialize our linear regression learner.

- Explain params to dump the parameters we can use.

```

import org.apache.spark.ml.regression.LinearRegression
import org.apache.spark.ml.regression.LinearRegressionModel
import org.apache.spark.ml.Pipeline
val lr = new LinearRegression()
lr.explainParams()

import org.apache.spark.ml.regression.LinearRegression
import org.apache.spark.ml.regression.LinearRegressionModel
import org.apache.spark.ml.Pipeline
lr: org.apache.spark.ml.regression.LinearRegression = linReg_301e5ead4215
res72: String =
elasticNetParam: the ElasticNet mixing parameter, in range [0, 1]. For alpha = 0, the penalty is an L2 penalty. For alpha = 1, it is an L1 penalty (default: 0.0)
featuresCol: features column name (default: features)
fitIntercept: whether to fit an intercept term (default: true)
labelCol: label column name (default: label)
maxIter: maximum number of iterations (>= 0) (default: 100)
predictionCol: prediction column name (default: prediction)
regParam: regularization parameter (>= 0) (default: 0.0)
solver: the solver algorithm for optimization. If this is not set or empty, default value is 'auto'. (default: auto)
standardization: whether to standardize the training features before fitting the model (default: true)
tol: the convergence tolerance for iterative algorithms (de...

```

#### 10. Set the parameters for the method:

- Set the name of the prediction column to "Predicted\_PE"
- Set the name of the label column to "PE"
- Set the maximum number of iterations to 100
- Set the regularization parameter to 0.1

```

lr.setPredictionCol("PE_Predict")
  .setLabelCol("PE")
  .setMaxIter(100)
  .setRegParam(0.1)

res74: org.apache.spark.ml.regression.LinearRegression = linReg_301e5ead4215

```

#### 11. Create a model by training on 'trainingSet'

- Use the new spark.ml pipeline API, Similar to python scikit-learn
- First train on the entire dataset to see what we get

```

val lrPipeline = new Pipeline()

```

```
lrPipeline.setStages(Array(vectorizer, lr))

val lrModel = lrPipeline.fit(trainingSet)
```

```
lrPipeline: org.apache.spark.ml.Pipeline = pipeline_5a953dde602d  
res76: lrPipeline.type = pipeline_5a953dde602d  
lrModel: org.apache.spark.ml.PipelineModel = pipeline_5a953dde602d
```

## 12. Linear Regression Model

Linear Regression is simply a Line of best fit over the data that minimizes the square of the error, given multiple input dimensions we can express each predictor as a line function of the form:

[ y = a + b x\_1 + b x\_2 + b x\_i ... \\ ]

where  $\backslash\backslash(a\backslash\backslash)$  is the intercept and the  $\backslash\backslash(b\backslash\backslash)$  are the coefficients.

```
val model = lrModel.stages(1).asInstanceOf[LinearRegressionModel]
```

```
model: org.apache.spark.ml.regression.LinearRegressionModel = linReg_30le5ead4215
```

13. To express the coefficients of that line we can retrieve the Estimator stage from the PipelineModel and express the weights and the intercept for the function.

- Intercepts

```
| val intercept = model.intercept
```

intercept: Double = 437.3755581407379

- Weights

```
val weights = model.weights.toArray
```

```
warning: there were 1 deprecation warning(s); re-run with -deprecation for details
weights: Array[Double] = Array(-1.9184240575231875, -0.2543479937050499, 0.07825010764537191, -0.1471027706567747)
```

#### 14. List of the column names (without PE)

```
| val featuresNoLabel = df.columns.filter(col => col != "PE")
```

```
featuresNoLabel: Array[String] = Array(AT, V, AP, RH)
```

#### 15. Merge the weights and labels

```
| val coefficents =
| sc.parallelize(weights).zip(sc.parallelize(featuresNoLabel))
| coefficents.collect.foreach(println)
```

```
coefficents: org.apache.spark.rdd.RDD[(Double, String)] = ZippedPartitionsRDD2[116] at zip at <console>:87
(-1.9184240575231875,AT)
(-0.2543479937050499,V)
(0.07825010764537191,AP)
(-0.1471027706567747,RH)
```

#### 16. Assigning value of intercept to equation

```
| var equation = s"y = $intercept "
```

```
equation: String = "y = 437.3755581407379 "
```

#### 17. Sort the coefficients from greatest absolute weight most to the least absolute weight

```
| var variables = Array
| coefficents.sortByKey().collect().foreach(x =>
| {
|   val weight = Math.abs(x._1)
|   val name = x._2
| }
```

```

    val symbol = if (x._1 > 0) "+" else "-"
    equation += (s" $symbol (${weight} * ${name})")
}
)

```

variables: Array.type = scala.Array@55b81d1f

## 18. Linear Regression Equation

```

    println("Linear Regression Equation: " + equation)

```

Linear Regression Equation:  $y = 437.3755581407379 - (1.9184240575231875 * AT) - (0.2543479937050499 * V) - (0.1471027706567747 * RH) + (0.07825010764537191 * AP)$

## 19. Apply our LR model to the test data and predict power output

```

    val predictionsAndLabels = lrModel.transform(testingSet)
    predictionsAndLabels.show(5)

```

Scroll through the resulting table and notice how the values in the Power Output (PE) column compare to the corresponding values in the predicted Power Output (Predicted\_PE) column.

```

predictionsAndLabels: org.apache.spark.sql.DataFrame = [AT: double, V: double, AP: double, RH: double, PE: double, features: vector, PE_Predict: double]
+---+---+---+---+-----+
| AT| V| AP| RH| PE|      features| PE_Predict|
+---+---+---+---+-----+
|3.38|41.31| 998.79|97.76|489.11|[3.38,41.31,998.7...| 484.1588273629686|
| 3.6|35.19|1018.73| 99.1|488.98|[3.6,35.19,1018.7...|486.65657322465705|
|3.73|39.42| 1024.4|82.42|488.58|[3.73,39.42,1024....|488.22863840871094|
|3.74|35.19|1018.58|98.84| 490.5|[3.74,35.19,1018....|486.41450306082777|
|3.95|38.44|1016.75|79.65|492.46|[3.95,38.44,1016....|487.86470750111897|
+---+---+---+---+-----+
only showing top 5 rows

```

## 20. Compute an evaluation metric for our test dataset

- Create an RMSE evaluator using the label and predicted columns
- Run the evaluator on the DataFrame
- Calculating R squared value.

**Note:** Another useful statistical evaluation metric is the coefficient of determination, denoted  $R^2$  or  $r^2$  and pronounced "R squared".

It is a number that indicates the proportion of the variance in the dependent variable that is predictable from the independent variable and it provides a measure of how well observed outcomes are replicated by the model, based on the proportion of total variation of outcomes explained by the model. The coefficient of determination ranges from 0 to 1 (closer to 1), and the higher the value, the better our model.

```
-----  
import org.apache.spark.mllib.evaluation.RegressionMetrics  
  
val metrics = new  
RegressionMetrics(predictionsAndLabels.select("PE_Predict",  
"PE").rdd.map(r => (r(0).asInstanceOf[Double],  
r(1).asInstanceOf[Double])))  
  
val rmse = metrics.rootMeanSquaredError  
  
val explainedVariance = metrics.explainedVariance  
  
val r2 = metrics.r2  
  
println (f"Root Mean Squared Error: $rmse")  
println (f"Explained Variance: $explainedVariance")  
println (f"R2: $r2")  
-----
```

```
-----  
import org.apache.spark.mllib.evaluation.RegressionMetrics  
metrics: org.apache.spark.mllib.evaluation.RegressionMetrics = org.apache.spark.mllib.evaluation.RegressionMetrics@f53cfef  
rmse: Double = 4.596932972254193  
explainedVariance: Double = 267.5580020580147  
r2: Double = 0.9275409236725611  
Root Mean Squared Error: 4.596932972254193  
Explained Variance: 267.5580020580147  
R2: 0.9275409236725611  
-----
```

**Note:** Generally, a good model will have 68% of predictions within 1 RMSE and 95% within 2 RMSE of the actual value. Let's calculate and see if a RMSE of 4.51 meets this criterion.

## Tuning and Evaluation

Now that we have a model with all of the data let's try to make a better model by tuning over several parameters to see if we can get better results.

```
-----  
import org.apache.spark.ml.tuning.{ParamGridBuilder, CrossValidator}  
import org.apache.spark.ml.evaluation.  
-----
```

```
import org.apache.spark.ml.tuning.{ParamGridBuilder, CrossValidator}
import org.apache.spark.ml.evaluation._
```

## 21. Use a cross validator to split the data into training and validation subsets

Let's set up our evaluator class to judge the model based on the best root mean squared error

```
val regEval = new RegressionEvaluator()
  regEval.setLabelCol("PE")
    .setPredictionCol("PE_Predict")
    .setMetricName("rmse")
```

```
regEval: org.apache.spark.ml.evaluation.RegressionEvaluator = regEval_6eb0d0675120
res98: regEval.type = regEval_6eb0d0675120
```

## 22. Create our crossvalidator with 3 fold cross validation

```
val crossval = new CrossValidator()
  crossval.setEstimator(lrPipeline)
  crossval.setNumFolds(5)
  crossval.setEvaluator(regEval)
```

```
crossval: org.apache.spark.ml.tuning.CrossValidator = cv_a61d2e94671b
res101: crossval.type = cv_a61d2e94671b
res102: crossval.type = cv_a61d2e94671b
res103: crossval.type = cv_a61d2e94671b
```

## 23. Tune over our regularization parameter from 0.01 to 0.10

```
val regParam = ((1 to 10) toArray).map(x => (x / 100.0))
```

```
warning: there were 1 feature warning(s); re-run with -feature for details
regParam: Array[Double] = Array(0.01, 0.02, 0.03, 0.04, 0.05, 0.06, 0.07, 0.08, 0.09, 0.1)
```

24. Create a parameter grid using the ParamGridBuilder, and add the grid to the CrossValidator.

```
val paramGrid = new ParamGridBuilder()
  .addGrid(lr.regParam, regParam)
  .build()

crossval.setEstimatorParamMaps(paramGrid)
```

```
paramGrid: Array[org.apache.spark.ml.param.ParamMap] =
Array({
  linReg_301e5ead4215-regParam: 0.01
}, {
  linReg_301e5ead4215-regParam: 0.02
}, {
  linReg_301e5ead4215-regParam: 0.03
}, {
  linReg_301e5ead4215-regParam: 0.04
}, {
  linReg_301e5ead4215-regParam: 0.05
}, {
  linReg_301e5ead4215-regParam: 0.06
}, {
  linReg_301e5ead4215-regParam: 0.07
}, {
  linReg_301e5ead4215-regParam: 0.08
},
```

25. Let's find and return the best model

```
val cvModel = crossval.fit(trainingSet)
```

```
cvModel: org.apache.spark.ml.tuning.CrossValidatorModel = cv_a61d2e94671b
```

26. We have tuned let's see what we got for tuning parameters and what our RMSE was versus our initial model.

- Use cvModel to compute an evaluation metric for our test dataset: testingSet

```
val predictionsAndLabels = cvModel.transform(testingSet)

val metrics = new
RegressionMetrics(predictionsAndLabels.select("PE_Predict",
```

```
"PE").rdd.map(r => (r(0).asInstanceOf[Double],  
r(1).asInstanceOf[Double]))
```

```
predictionsAndLabels: org.apache.spark.sql.DataFrame = [AT: double, V: double, AP: double, RH: double, PE: double, features: vector, PE_Predict: double]  
metrics: org.apache.spark.mllib.evaluation.RegressionMetrics = org.apache.spark.mllib.evaluation.RegressionMetrics@1a3b890e
```

## 27. Run the previously created RMSE evaluator, on the predictionsAndLabels DataFrame

- compute the r2 evaluation metric for our test dataset

```
val rmse = metrics.rootMeanSquaredError  
val explainedVariance = metrics.explainedVariance  
val r2 = metrics.r2  
  
println (f"Root Mean Squared Error: $rmse")  
println (f"Explained Variance: $explainedVariance")  
println (f"R2: $r2")
```

```
rmse: Double = 4.593657981825083  
explainedVariance: Double = 269.9212285609958  
r2: Double = 0.9276441308553334  
Root Mean Squared Error: 4.593657981825083  
Explained Variance: 269.9212285609958  
R2: 0.9276441308553334
```

So our initial untuned and tuned linear regression models are statistically identical.

## Regression Tree

Given that the only linearly correlated variable is Temperature, it makes sense try another machine learning method such a Decision Tree to handle non-linear data and see if we can improve our model

A Decision Tree creates a model based on splitting variables using a tree structure. We will first start with a single decision tree model.

## 1. Create a DecisionTreeRegressor

```
import org.apache.spark.ml.regression.DecisionTreeRegressor

val dt = new DecisionTreeRegressor()
dt.setLabelCol("PE")
dt.setPredictionCol("PE_Predict")
dt.setFeaturesCol("features")
dt.setMaxBins(100)
```

```
import org.apache.spark.ml.regression.DecisionTreeRegressor
dt: org.apache.spark.ml.regression.DecisionTreeRegressor = dtr_8585361dd982
res117: org.apache.spark.ml.regression.DecisionTreeRegressor = dtr_8585361dd982
res118: org.apache.spark.ml.regression.DecisionTreeRegressor = dtr_8585361dd982
res119: org.apache.spark.ml.regression.DecisionTreeRegressor = dtr_8585361dd982
res120: dt.type = dtr_8585361dd982
```

## 2. Create a Pipeline and set the stages of the pipeline

```
val dtPipeline = new Pipeline()
dtPipeline.setStages(Array(vectorizer, dt))
```

## 3. Let's reuse our CrossValidator

- Use [ParamGridBuilder] to build a parameter grid with the parameter `dt.maxDepth` and a list of the values 2 and 3, and add the grid to the [CrossValidator].
- Run the [CrossValidator] to find the parameters that yield the best model (i.e. lowest RMSE) and return the best model.

```
crossval.setEstimator(dtPipeline)

val paramGrid = new ParamGridBuilder()
.addGrid(dt.maxDepth, Array(2, 3))
.build()

crossval.setEstimatorParamMaps(paramGrid)

val dtModel = crossval.fit(trainingSet)
```

```

dtPipeline: org.apache.spark.ml.Pipeline = pipeline_78ed1072fa9d
res122: dtPipeline.type = pipeline_78ed1072fa9d
res123: crossval.type = cv_a61d2e94671b
paramGrid: Array[org.apache.spark.ml.param.ParamMap] =
  Array({
    dtr_8585361dd982-maxDepth: 2
  }, {
    dtr_8585361dd982-maxDepth: 3
  })
res124: crossval.type = cv_a61d2e94671b
dtModel: org.apache.spark.ml.tuning.CrossValidatorModel = cv_a61d2e94671b

```

Let's see how our tuned DecisionTreeRegressor model's RMSE and  $\sqrt{r^2}$  values compare to our tuned LinearRegression model.

```

import org.apache.spark.ml.regression.DecisionTreeRegressionModel
import org.apache.spark.ml.PipelineModel

val predictionsAndLabels = dtModel.bestModel.transform(testingSet)

val metrics = new
RegressionMetrics(predictionsAndLabels.select("PE_Predict",
"PE").map(r => (r(0).asInstanceOf[Double],
r(1).asInstanceOf[Double])))

val rmse = metrics.rootMeanSquaredError
val explainedVariance = metrics.explainedVariance
val r2 = metrics.r2

println (f"Root Mean Squared Error: $rmse")
println (f"Explained Variance: $explainedVariance")
println (f"R2: $r2")

```

```

import org.apache.spark.ml.regression.DecisionTreeRegressionModel
import org.apache.spark.ml.PipelineModel
predictionsAndLabels: org.apache.spark.sql.DataFrame = [AT: double, V: double, AP: double, RH: double, PE: double, features: vector, PE_Predict: double]
metrics: org.apache.spark.mllib.evaluation.RegressionMetrics = org.apache.spark.mllib.evaluation.RegressionMetrics@54e2ea3
rmse: Double = 5.169951561569703
explainedVariance: Double = 264.7744220845704
r2: Double = 0.9083506479156687
Root Mean Squared Error: 5.169951561569703
Explained Variance: 264.7744220845704
R2: 0.9083506479156687

```

This line will pull the Decision Tree model from the Pipeline as display it as an if-then-else string

```
| dtModel.bestModel.asInstanceOf[PipelineModel].stages.last.asInstanceOf[DecisionTreeRegressionModel].toDebugString
```

```
res132: String =
"DecisionTreeRegressionModel (uid=dtr_8585361dd982) of depth 3 with 15 nodes
 If (feature 0 <= 17.85)
   If (feature 0 <= 11.67)
     If (feature 0 <= 8.63)
       Predict: 483.79243774574053
     Else (feature 0 > 8.63)
       Predict: 476.3108737113401
   Else (feature 0 > 11.67)
     If (feature 0 <= 14.35)
       Predict: 469.0458953626635
     Else (feature 0 > 14.35)
       Predict: 462.1050880962613
   Else (feature 0 > 17.85)
     If (feature 0 <= 22.99)
       If (feature 1 <= 45.87)
         Predict: 458.12494773519165
```

So our DecisionTree has slightly worse RMSE than our LinearRegression model (LR: 4.59 vs DT: 5.19).

## Exercise9: Machine Learning

To perform below exercise, install R-studio from ‘software’ folder or open R console and run below commands from command line.

Alternatively, while dealing with huge datasets we will use ‘Zeppelin’ with R interpreter pre-installed on the machine.

Statistical analysis in R is performed by using many in-built functions.

Functions are part of the R base package. These functions take R vector as an input along with the arguments and give the result.

### 1. Mean

From R console, use function mean() to calculate mean.

```

# Create a vector.
x <- c(12,7,3,4.2,18,2,54,-21,8,-5)
# Find Mean.
result.mean <- mean(x)
print(result.mean)

> x <- c(12,7,3,4.2,18,2,54,-21,8,-5)
> result.mean <- mean(x)
> print(result.mean)
[1] 8.22

```

## 2. Median

Similarly, to calculate median

```

# Create the vector.
x <- c(12,7,3,4.2,18,2,54,-21,8,-5)
# Find the median.
median.result <- median(x)
print(median.result)

> median.result <- median(x)
> print(median.result)
[1] 5.6
>

```

## Linear Regression

Exercise9.1- Using R-studio/ R-console

**Example:** Predicting the weight of a person when his height is known.

**Steps:**

1. Gather the sample of observed values of height and corresponding weight.

**Input:**

```

# Values of height
151, 174, 138, 186, 128, 136, 179, 163, 152, 131

```

```
| # Values of weight.  
|  
| 63, 81, 56, 91, 47, 57, 76, 72, 62, 48
```

## 2. Create relationship model using **lm()** function.

Find the coefficients from the model and create the equation.

```
| x <- c(151, 174, 138, 186, 128, 136, 179, 163, 152, 131)  
|  
| y <- c(63, 81, 56, 91, 47, 57, 76, 72, 62, 48)  
|  
| # Apply the lm() function.  
|  
| relation <- lm(y~x)  
|  
| print(relation)
```

```
> x <- c(151, 174, 138, 186, 128, 136, 179, 163, 152, 131)  
> y <- c(63, 81, 56, 91, 47, 57, 76, 72, 62, 48)  
> relation <- lm(y~x)  
> print(relation)  
  
Call:  
lm(formula = y ~ x)  
  
Coefficients:  
(Intercept)           x  
-38.4551        0.6746
```

## 3. Get summary of the model to know average error i.e. **residuals**.

```
| x <- c(151, 174, 138, 186, 128, 136, 179, 163, 152, 131)  
|  
| y <- c(63, 81, 56, 91, 47, 57, 76, 72, 62, 48)  
|  
| # Apply the lm() function.  
|  
| relation <- lm(y~x)  
|  
| print(summary(relation))
```

```

> relation <- lm(y~x)
> print(summary(relation))

Call:
lm(formula = y ~ x)

Residuals:
    Min      1Q  Median      3Q     Max 
-6.3002 -1.6629  0.0412  1.8944  3.9775 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) -38.45509   8.04901 -4.778  0.00139 ***
x             0.67461   0.05191 12.997 1.16e-06 ***
---
Signif. codes:  0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Residual standard error: 3.253 on 8 degrees of freedom
Multiple R-squared:  0.9548, Adjusted R-squared:  0.9491 
F-statistic: 168.9 on 1 and 8 DF,  p-value: 1.164e-06

```

4. To predict the weight of new person use **predict()** function.

```

# The predictor vector.

x <- c(151, 174, 138, 186, 128, 136, 179, 163, 152, 131)

# The response vector.

y <- c(63, 81, 56, 91, 47, 57, 76, 72, 62, 48)

# Apply the lm() function.

relation <- lm(y~x)

# Find weight of a person with height 170.

a <- data.frame(x = 170)

result <- predict(relation,a)

print(result)

```

```

> relation <- lm(y~x)
> a <- data.frame(x = 170)
> result <- predict(relation,a) .
> print(result)
1
76.22869
> |

```

5. Visualize the regression graphically

```

# Create the predictor and response variable.

x <- c(151, 174, 138, 186, 128, 136, 179, 163, 152, 131)

```

```

y <- c(63, 81, 56, 91, 47, 57, 76, 72, 62, 48)
relation <- lm(y~x)

# Give the chart file a name.
png(file = "linearregression.png")

# Plot the chart.
plot(y,x,col = "blue",main = "Height & Weight Regression",
abline(lm(x~y)),cex = 1.3,pch = 16,xlab = "Weight in Kg",ylab =
"Height in cm")

# Save the file.

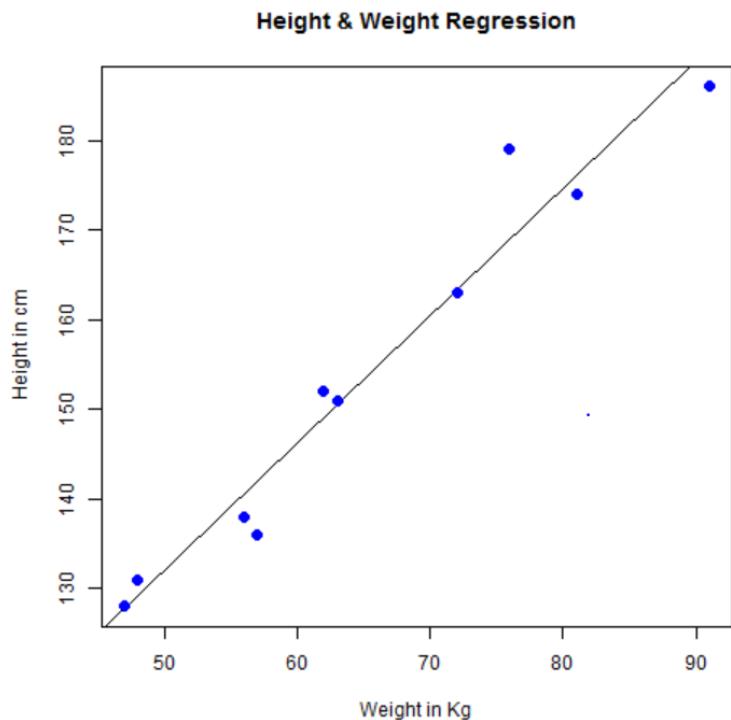
dev.off()

```

```

> relation <- lm(y~x)
> png(file = "linearregression.png")
> plot(y,x,col = "blue",main = "Height & Weight Regression",abline(lm(x~y)),cex = 1.3,pch = 16,xlab
= "Weight in Kg",ylab = "Height in cm")
> dev.off()
null device
      1
>

```



## Exercise9.2 Using zeppelin

### Goal:

1. Start with importing the required package for R

```
%r require(ggplot2)
```

2. Load the data and list the features

```
%r

power_plant = read.table("/data/CCPP/Folds5x2_pp_1.csv", header =
TRUE, sep = ",")  
str(power_plant)

# Features consist of hourly average ambient variables  
# Temperature (T) in the range 1.81°C and 37.11°C,  
# Ambient Pressure (AP) in the range 992.89-1033.30 milibar,  
# Relative Humidity (RH) in the range 25.56% to 100.16%  
# Exhaust Vacuum (V) in the range 25.36-81.56 cm Hg  
# Net hourly electrical energy output (EP) 420.26-495.76 MW
```

```
'data.frame': 9568 obs. of 5 variables:  
 $ AT: num 14.96 25.18 5.11 20.86 10.82 ...  
 $ V : num 41.8 63 39.4 57.3 37.5 ...  
 $ AP: num 1024 1020 1012 1010 1009 ...  
 $ RH: num 73.2 59.1 92.1 76.6 96.6 ...  
 $ PE: num 463 444 489 446 474 ...
```

3. To check first few entries

```
%r head(power_plant)
```

```
1 14.96 41.76 1024.07 73.17 463.26
2 25.18 62.96 1020.04 59.08 444.37
3 5.11 39.40 1012.16 92.14 488.56
4 20.86 57.32 1010.24 76.64 446.48
5 10.82 37.50 1009.23 96.62 473.90
6 26.27 59.44 1012.23 58.77 443.67
```

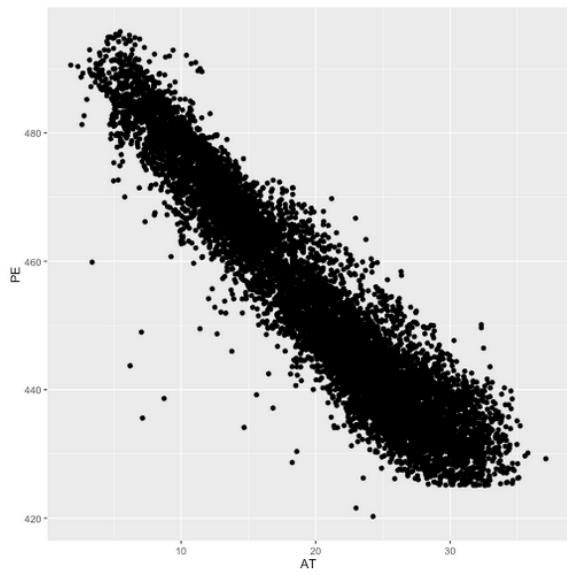
#### 4. To check the summary of the dataset

```
%r summary(power_plant)
```

```
Min.   : 1.81   Min.   :25.36   Min.   : 992.9   Min.   : 25.56<br />
1st Qu.:13.51   1st Qu.:41.74   1st Qu.:1009.1   1st Qu.: 63.33<br />
Median :20.34   Median :52.08   Median :1012.9   Median : 74.97<br />
Mean   :19.65   Mean   :54.31   Mean   :1013.3   Mean   : 73.31<br />
3rd Qu.:25.72   3rd Qu.:66.54   3rd Qu.:1017.3   3rd Qu.: 84.83<br />
Max.   :37.11   Max.   :81.56   Max.   :1033.3   Max.   :100.16<br />
PE<br />
Min.   :420.3<br />
1st Qu.:439.8<br />
Median :451.6<br />
Mean   :454.4<br />
3rd Qu.:468.4<br />
Max.   :495.8
```

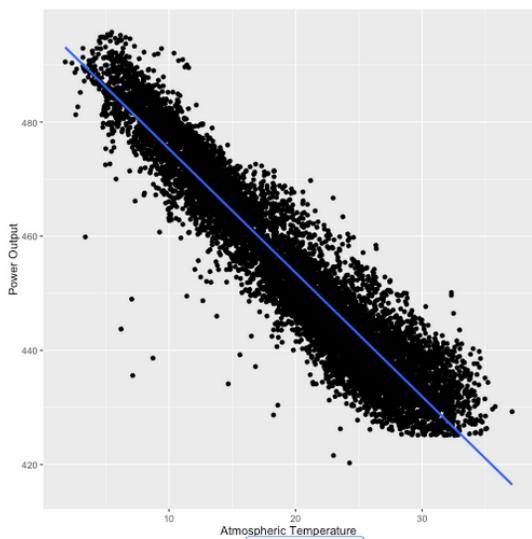
#### 5. Plot between AT and PE

```
%r
ggplot(power_plant, aes(x = AT, y = PE)) +
  geom_point()
```



## 6. Another way of representation

```
%r  
ggplot(power_plant, aes(x = AT, y = PE)) +  
  geom_point() +  
  geom_smooth(method = "lm") +  
  labs(x = "Atmospheric Temperature", y = "Power Output")
```



## Linear Regression Model

```
%r  
  
lm.fit = lm(PE ~ AT, data = power_plant)  
#Statistical Significance of the Model  
capture.output(summary(lm.fit))  
  
[2] "Call:<br />  
[3] "lm(formula = PE ~ AT, data = power_plant)"<br />  
[4] "<br />  
[5] "Residuals:<br />  
[6] "   Min     1Q  Median      3Q     Max "<br />  
[7] "-45.951 -3.644   0.101   3.696  23.251 "<br />  
[8] "<br />  
[9] "Coefficients:<br />  
[10] "             Estimate Std. Error t value Pr(>|t|)    "<br />  
[11] "(Intercept) 497.034120  0.156434 3177.3  &lt;2e-16 <em>***<br />  
[12] "AT         -2.171320  0.007443  -291.7  &lt;2e-16 *</em><em>"<br />  
[13] "-<br />  
[14] "Signif. codes:  0 '(<em><strong> 0.001 '</strong>' 0.01 '*' 0.05 '.' 0.1 ' ' 1"  
[15] "<br />  
[16] "Residual standard error: 5.426 on 9566 degrees of freedom"<br />  
[17] "Multiple R-squared:  0.8989, \tAdjusted R-squared:  0.8989 "<br />  
[18] "F-statistic: 8.51e+04 on 1 and 9566 DF,  p-value: &lt; 2.2e-16"<br />  
[19] "
```



## Interpretation

### Residuals:

The residuals are the difference between the actual values of the variable you're predicting and predicted values from your regression— $y - \hat{y}$ . For most regressions you want your residuals to look like a normal distribution when plotted. If our residuals are normally distributed, this indicates the mean of the difference between our predictions and the actual values is close to 0 (good) and that when we miss, we're missing both short and long of the actual value, and the likelihood of a miss being far from the actual value gets smaller as the distance from the actual value gets larger.

### Significance Stars

The stars are shorthand for significance levels, with the number of asterisks displayed according to the p-value computed. for high significance and \* for low significance.

### Estimated Coefficient

The estimated coefficient is the value of slope calculated by the regression. It might seem a little confusing that the Intercept also has a value, but just think of it as a slope that is always multiplied by 1.

### **Standard Error of the Coefficient Estimate**

Measure of the variability in the estimate for the coefficient. Lower means better but this number is relative to the value of the coefficient. As a rule of thumb, you'd like this value to be at least an order of magnitude less than the coefficient estimate.

### **t-value or t-statistics**

Score that measures whether or not the coefficient for this variable is meaningful for the model. You probably won't use this value itself, but know that it is used to calculate the p-value and the significance levels.

### **Significance Legend**

The more punctuation there is next to your variables, the better.

Blank=bad,

Dots=pretty good

Stars=good

More Stars=very good

### **Residual Std Error / Degrees of Freedom**

The Residual Std Error is just the standard deviation of your residuals. You'd like this number to be proportional to the quantiles of the residuals in #1. For a normal distribution, the 1st and 3rd quantiles should be  $1.5 \pm$  the std error.

The Degrees of Freedom is the difference between the number of observations included in your training sample and the number of variables used in your model (intercept counts as a variable).

### **Variable p-value**

Probability the variable is NOT relevant. You want this number to be as small as possible. If the number is really small, R will display it in scientific notation.

### **R-squared**

Metric for evaluating the goodness of fit of your model. Higher is better with 1 being the best. An R<sup>2</sup> value close to 1 indicates that the model explains a large portion of the variance in the response variable.

WARNING: While a high R-squared indicates good correlation, correlation does not always imply causation.

### **R-square adjusted**

R-square adjusted is penalized for having a large number of parameters in the model

$$R^2_{\text{adjusted}} = 1 - \frac{(1 - R^2)(n - 1)}{(n - p)}$$

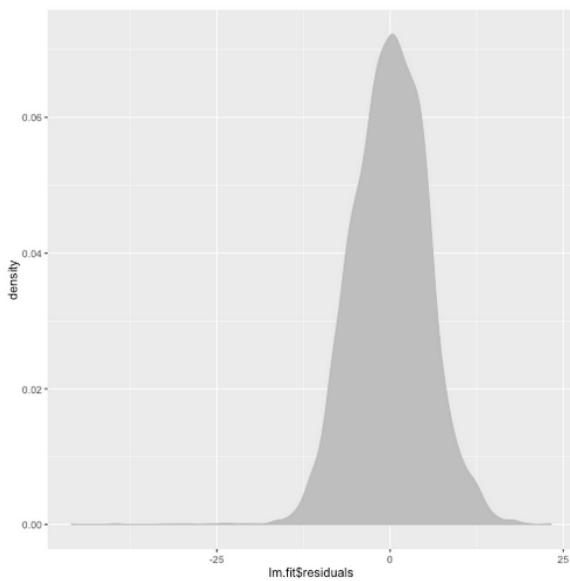
### **F-statistic & resulting p-value**

$$F = (\text{explained variance}) / (\text{unexplained variance})$$

Performs an F-test on the model. This takes the parameters of our model (in our case we only have 1) and compares it to a model that has fewer parameters. In theory the model with more parameters should fit better. If the model with more parameters (your model) doesn't perform better than the model with fewer parameters, the F-test will have a high p-value (probability NOT significant boost). If the model with more parameters is better than the model with fewer parameters, you will have a lower p-value.

When there is no relationship between the response and predictors, one would expect the F-statistic to take on a value close to 1. Large F-statistic suggests that at least some predictors are related to the response variable.

```
%r ggplot() + geom_density(aes(lm.fit$residuals), fill = "grey", color = "grey")
```



```
%r #Residual standard error:
```

```
sd(lm.fit$residuals)
```

```
[1] 5.425383
```

```
%r #R2:
```

```
R2 = sum((lm.fit$fitted.values - mean(power_plant$PE)) ^ 2) /  
sum((power_plant$PE - mean(power_plant$PE)) ^ 2)
```

```
R2
```

```
[1] 0.8989476
```

```
%r #F-statistics
```

```
# F = ((TSS - RSS) / p) / (RSS / (n - p - 1))
```

```
# TSS = Total sum of square = sum((y - mean(y)) ^ 2)
```

```

# RSS = Residual Sum of Square = sum((y_predicted - y) ^ 2)
# p = no of predictors
# n = number of observations

p = 1
TSS = sum((power_plant$PE - mean(power_plant$PE)) ^ 2)
RSS = sum((lm.fit$fitted.values - power_plant$PE) ^ 2)

n = 9568
F = ((TSS - RSS) / p) / (RSS / (n - p - 1))
F

```

[1] 85097.76

```

%r #R-squared Adjusted

R2.adjusted = 1 - (1 - R2)*(n - 1) / (n - p)
R2.adjusted

```

[1] 0.8989476

```

%r names(lm.fit)

[1] "coefficients"   "residuals"      "effects"        "rank"
[5] "fitted.values"  "assign"         "qr"             "df.residual"
[9] "xlevels"         "call"          "terms"          "model"

```

```
%r coef(lm.fit)
```

	(Intercept)	AT
	497.03412	-2.17132

```
%r confint(lm.fit)
```

```
(Intercept) 496.72748 497.34076  
AT          -2.18591  -2.15673
```

## Logistic Regression

The Logistic Regression is a regression model in which the response variable (dependent variable) has categorical values such as True/False or 0/1.

### **Example:**

In "mtcars" data set, the transmission mode (automatic or manual) is described by the column am which is a binary value (0 or 1). We can create a logistic regression model between the columns "am" and 3 other columns - hp, wt and cyl.

```
# Select some columns form mtcars.  
  
input <- mtcars[,c("am","cyl","hp","wt")]  
  
print(head(input))  
  
> input <- mtcars[,c("am","cyl","hp","wt")]  
> print(head(input))  
      am cyl  hp   wt  
Mazda RX4       1   6 110 2.620  
Mazda RX4 Wag   1   6 110 2.875  
Datsun 710      1   4  93 2.320  
Hornet 4 Drive   0   6 110 3.215  
Hornet Sportabout 0   8 175 3.440  
Valiant        0   6 105 3.460
```

### Create Regression Model:

```
input <- mtcars[,c("am","cyl","hp","wt")]  
  
am.data = glm(formula = am ~ cyl + hp + wt, data = input, family = binomial)
```

```

print(summary(am.data))

> input <- mtcars[,c("am","cyl","hp","wt")]
> am.data = glm(formula = am ~ cyl + hp + wt, data = input, family = binomial)
> print(summary(am.data))

Call:
glm(formula = am ~ cyl + hp + wt, family = binomial, data = input)

Deviance Residuals:
    Min      1Q  Median      3Q     Max 
-2.17272 -0.14907 -0.01464  0.14116  1.27641 

Coefficients:
            Estimate Std. Error z value Pr(>|z|)    
(Intercept) 19.70288   8.11637  2.428   0.0152 *  
cyl          0.48760   1.07162  0.455   0.6491    
hp           0.03259   0.01886  1.728   0.0840 .  
wt          -9.14947   4.15332 -2.203   0.0276 *  
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 43.2297 on 31 degrees of freedom
Residual deviance: 9.8415 on 28 degrees of freedom
AIC: 17.841

Number of Fisher Scoring iterations: 8

```

## Conclusion:

As the p-value in the last column is more than 0.05 for the variables "cyl" and "hp", we consider them to be insignificant in contributing to the value of the variable "am". Only weight (wt) impacts the "am" value in this regression model.

## Decision Tree

Decision tree is a graph to represent choices and their results in form of a tree. The nodes in the graph represent an event or choice and the edges of the graph represent the decision rules or conditions.

The R package "**party**" is used to create decision trees.

```

install.packages("party")

```

## Input Data

We will use the R in-built data set named **readingSkills** to create a decision tree. It describes the score of someone's readingSkills if we know the variables "age", "shoe size", "score" and whether the person is a native speaker or not.

```
| # Load the party package. It will automatically load other dependent
| packages.
|
| library(party)
```

```
| # Print some records from data set readingSkills.
|
| print(head(readingSkills))
```

```
> library(party)
Loading required package: grid
Loading required package: mvtnorm
Loading required package: modeltools
Loading required package: stats4
Loading required package: strucchange
Loading required package: zoo
```

```
Attaching package: 'zoo'
```

```
The following objects are masked from 'package:base':
```

```
  as.Date, as.Date.numeric
```

```
>Loading required package: sandwich
> print(head(readingSkills))
  nativeSpeaker age shoeSize   score
1       yes     5 24.83189 32.29385
2       yes     6 25.95238 36.63105
3       no      11 30.42170 49.60593
4      yes     7 28.66450 40.28456
5      yes    11 31.88207 55.46085
6      yes    10 30.07843 52.83124
```

We will use the **ctree()** function to create the decision tree and see its graph.

```
| # Load the party package. It will automatically load other dependent
| packages.
```

```
| library(party)
```

```
| # Create the input data frame.
```

```
| input.dat <- readingSkills[c(1:105),]
```

```

# Give the chart file a name.

png(file = "decision_tree.png")

# Create the tree.

output.tree <- ctree(
  nativeSpeaker ~ age + shoeSize + score,
  data = input.dat)

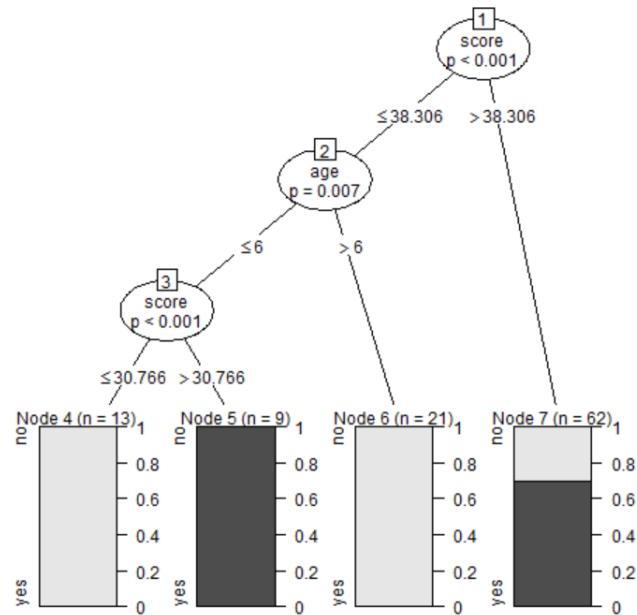
# Plot the tree.

plot(output.tree)

# Save the file.

dev.off()

```



## Conclusion

From the decision tree shown above we can conclude that anyone whose readingSkills score is less than 38.3 and age is more than 6 is not a native Speaker.

## Random Forest

In the random forest approach, a large number of decision trees are created. Every observation is fed into every decision tree. The most common outcome for each observation is used as the final output. A new observation is fed into all the trees and taking a majority vote for each classification model.

An error estimate is made for the cases which were not used while building the tree. That is called an **OOB (Out-of-bag)** error estimate which is mentioned as a percentage.

The R package "**randomForest**" is used to create random forests.

```
install.packages("randomForest")  
  
> install.packages("randomForest")  
trying URL 'https://cran.rstudio.com/bin/windows/contrib/3.3/randomForest_4.6-12.zip'  
Content type 'application/zip' length 177129 bytes (172 KB)  
downloaded 172 KB  
  
package 'randomForest' successfully unpacked and MD5 sums checked  
  
The downloaded binary packages are in  
C:\Users\technocracy\AppData\Local\Temp\RtmpcJlGBQ\downloaded_packages
```

## Input Data

We will use the R in-built data set named `readingSkills` to create a decision tree. It describes the score of someone's `readingSkills` if we know the variables "age", "shoesize", "score" and whether the person is a native speaker.

```
# Load the party package. It will automatically load other required  
packages.  
  
library(party)  
  
# Print some records from data set readingSkills.  
  
print(head(readingSkills))
```

We will use the **randomForest()** function to create the decision tree and see it's graph.

```
# Load the party package. It will automatically load other required
packages.

library(party)

library(randomForest)

# Create the forest.

output.forest <- randomForest(nativeSpeaker ~ age + shoeSize + score,
                                data = readingSkills)

# View the forest results.

print(output.forest)

# Importance of each predictor.

print(importance(fit,type = 2))

> print(output.forest)

Call:
randomForest(formula = nativeSpeaker ~ age + shoeSize + score,      data = readingSkills)
                         Type of random forest: classification
                               Number of trees: 500
No. of variables tried at each split: 1

          OOB estimate of  error rate: 1%
Confusion matrix:
  no yes class.error
no  99   1       0.01
yes   1  99       0.01
> |
```

## Conclusion

From the random forest shown above we can conclude that the shoesize and score are the important factors deciding if someone is a native speaker or not. Also the model has only 1% error which means we can predict with 99% accuracy.

## Classification

### **About Dataset: D2hawkEye**

D2Hawkeye, Inc. develops data service solutions for the healthcare industry.

D2ReportManager that helps to design, distribute, and archive financial and clinical reports; D2BenefitAdvisor and D2Analyzer, on-line analytical processing tools that enable users to build data sets for data mining and ad hoc query analysis; D2Consulting, a medical and clinical advisory service; and D2Connect that provides e-business services, custom software development, and technology outsourcing. It also provides medical advisory, medical analytics, benchmarking, and other business services on an outsourced basis. The company serves health plans, third party administrators, case management and disease management industry, benefit consultants, physician groups, employers, and pharmaceutical industry, as well as re-insurers, MGUs, and stop-loss groups.

### Mission

D2Hawkeye tries to improve healthcare case management

- Identify high-risk patients
- Work with patients to manage treatment and associated costs
- Arrange specialist care
- Medical costs often relate to severity of health problems, and are an issue for both patient and provider

**Goal:** Improve the quality of cost predictions

1. Load the data from local filesystem.

```
| %r Claims = read.csv("data/ClaimsData.csv")|
```

2. List the variables

```
| %r str(Claims)|
```

```
'data.frame': 458005 obs. of 16 variables:
 $ age          : int 85 59 67 52 67 68 75 70 67 67 ...
 $ alzheimers   : int 0 0 0 0 0 0 0 0 0 0 ...
 $ arthritis    : int 0 0 0 0 0 0 0 0 0 0 ...
 $ cancer       : int 0 0 0 0 0 0 0 0 0 0 ...
 $ copd          : int 0 0 0 0 0 0 0 0 0 0 ...
 $ depression   : int 0 0 0 0 0 0 0 0 0 0 ...
 $ diabetes     : int 0 0 0 0 0 0 0 0 0 0 ...
 $ heart.failure: int 0 0 0 0 0 0 0 0 0 0 ...
 $ ihd           : int 0 0 0 0 0 0 0 0 0 0 ...
 $ kidney        : int 0 0 0 0 0 0 0 0 0 0 ...
 $ osteoporosis  : int 0 0 0 0 0 0 0 0 0 0 ...
 $ stroke        : int 0 0 0 0 0 0 0 0 0 0 ...
 $ reimbursement2008: int 0 0 0 0 0 0 0 0 0 0 ...
 $ bucket2008   : int 1 1 1 1 1 1 1 1 1 1 ...
 $ reimbursement2009: int 0 0 0 0 0 0 0 0 0 0 ...
 $ bucket2009   : int 1 1 1 1 1 1 1 1 1 1 ...
```

### 3. List the output of \$bucket2009

```
%r capture.output(base:::table(Claims$bucket2009))
```

```
[1] "<br />
[2] " 1   2   3   4   5 "
[3] "307444 87099 40976 19843 2643 "
```

Took a few seconds. Last updated by anonymous at September 06 2016, 10:32:33 PM.

### 4. View the summary of claims dataset

```
%r summary(Claims)
```

```
Min.   : 26.00  Min.   :0.0000  Min.   :0.0000  Min.   :0.0000<br />
1st Qu.: 67.00  1st Qu.:0.0000  1st Qu.:0.0000  1st Qu.:0.0000<br />
Median : 73.00  Median :0.0000  Median :0.0000  Median :0.0000<br />
Mean   : 72.63  Mean   :0.1922  Mean   :0.1543  Mean   :0.06411<br />
3rd Qu.: 81.00  3rd Qu.:0.0000  3rd Qu.:0.0000  3rd Qu.:0.0000<br />
Max.   :100.00  Max.   :1.0000  Max.   :1.0000  Max.   :1.0000<br />
      copd      depression      diabetes      heart.failure<br />
Min.   :0.0000  Min.   :0.0000  Min.   :0.0000  Min.   :0.0000<br />
1st Qu.:0.0000 1st Qu.:0.0000  1st Qu.:0.0000  1st Qu.:0.0000<br />
Median :0.0000  Median :0.0000  Median :0.0000  Median :0.0000<br />
Mean   :0.1361  Mean   :0.2131  Mean   :0.3805  Mean   :0.2847<br />
3rd Qu.:0.0000  3rd Qu.:0.0000  3rd Qu.:1.0000  3rd Qu.:1.0000<br />
Max.   :1.0000  Max.   :1.0000  Max.   :1.0000  Max.   :1.0000<br />
      ihd      kidney      osteoporosis      stroke<br />
Min.   :0.0000  Min.   :0.0000  Min.   :0.0000  Min.   :0.0000<br />
1st Qu.:0.0000 1st Qu.:0.0000  1st Qu.:0.0000  1st Qu.:0.0000<br />
Median :0.0000  Median :0.0000  Median :0.0000  Median :0.0000<br />
      ..      ..      ..      ..      ..      ..      ..      ..
```

### 5. Percentage of patients in each cost bucket

```
%r base:::table(Claims$bucket2009)/nrow(Claims)
```

```
0.671267781 0.190170413 0.089466272 0.043324855 0.005770679
```

### 6. Model for healthcare cost prediction

**Goal:** To predict the cost bucket the patient fell into in 2009 using the CART model.

```
%r require(caTools) #For splitting data into training and test sets
```

7. Split the entire data into training and test sets, with 70/30 split.

```
%r  
set.seed(88)  
split = sample.split(Claims$bucket2009, SplitRatio = 0.7)  
ClaimsTrain = subset(Claims, split == TRUE)  
ClaimsTest = subset(Claims, split == FALSE)  
Summary(ClaimsTrain)
```

What is the average age of patients in the training set, ClaimsTrain?

```
mean(ClaimsTrain$age)
```

### Baseline Method and Penalty Matrix

The baseline method would predict that the cost bucket for a patient in 2009 will be the same as it was in 2008.

So let's create a classification matrix to compute the accuracy for the baseline method on the test set.

The accuracy is the sum of the diagonal, the observations that were classified correctly, divided by the total number of observations in our test set.

```
%r  
confMatrix = base::table(ClaimsTest$bucket2009, ClaimsTest$bucket2008)  
confMatrix
```

```
1 110138  7787  3427  1452   174  
2 16000   10721  4629  2931   559  
3 7006    4629  2774  1621   360  
4 2688    1943  1415  1539   352  
5 293     191   160   369   104
```

```
%r diag(confMatrix)
```

```
110138 10721 2774 1539 104
```

8. Baseline accuracy on test set

```
%r  
sum(diag(confMatrix)) / nrow(ClaimsTest)
```

```
[1] 0.6838135
```

## 9. Define a “penaltymatrix” as the cost of being wrong

```
%r  
PenaltyMatrix =  
matrix(c(0,1,2,3,4,2,0,1,2,3,4,2,0,1,2,6,4,2,0,1,8,6,4,2,0), byrow = TRUE,  
nrow = 5)  
PenaltyMatrix  
# X axis corresponds to the predictions and Y axis to the actuals
```

```
[1,]   0   1   2   3   4  
[2,]   2   0   1   2   3  
[3,]   4   2   0   1   2  
[4,]   6   4   2   0   1  
[5,]   8   6   4   2   0
```

## 10. Penalty error of Baseline Method

```
%r  
penaltyErrors = as.matrix(base::table(ClaimsTest$bucket2009,  
ClaimsTest$bucket2008)) * PenaltyMatrix  
penaltyErrors
```

```
1   0 7787 6854 4356 696  
2 32000   0 4629 5862 1677  
3 28024 9258   0 1621 720  
4 16128 7772 2830   0 352  
5 2344 1146 640  618   0
```

```
%r #Penalty Error Rate based on baseline data  
sum(penaltyErrors) / nrow(ClaimsTest)
```

```
[1] 0.7386055
```

## Our Goals:

Now our goal is to create a CART model that will give better accuracy higher than 68.38% and a penalty error lower than 0.739

```
%r
```

```
# Load necessary libraries  
library(rpart) #For building cart model  
library(rpart.plot) #For plotting cart model
```

### Complexity Parameter using Cross Validation

The complexity parameter (cp) is used to control the size of the decision tree and to select the optimal tree size.

```
%r  
  
library(e1071)  
library(caret)  
  
%r  
  
numFolds = trainControl( method = "cv", number = 10 )  
cpGrid = expand.grid( .cp = seq(0.0,0.03,0.001))
```

### 11. Perform the cross validation

```
cv = train(bucket2009 ~ age + alzheimers + arthritis + cancer + copd +  
depression + diabetes + heart.failure + ihd + kidney + osteoporosis +  
stroke + bucket2008 + reimbursement2008, data = ClaimsTrain, method =  
"rpart", trControl = numFolds, tuneGrid = cpGrid )
```

```
%r  
  
capture.output(cv)
```

```
[1] "CART "<br />
[2] "<br />
[3] "274803 samples"<br />
[4] "    15 predictor"<br />
[5] "<br />
[6] "No pre-processing"<br />
[7] "Resampling: Cross-Validated (10 fold) "<br />
[8] "Summary of sample sizes: 247322, 247322, 247323, 247322, 247322, 247323, ... "
[9] "Resampling results across tuning parameters:"<br />
[10] "<br />
[11] "   cp      RMSE      Rsquared "<br />
[12] " 0.000  0.7938988  0.2168627"<br />
[13] " 0.001  0.7276304  0.3012112"<br />
[14] " 0.002  0.7295325  0.2975572"<br />
[15] " 0.003  0.7347164  0.2875466"<br />
[16] " 0.004  0.7364503  0.2841780"<br />
[17] " 0.005  0.7388593  0.2794822"<br />
```

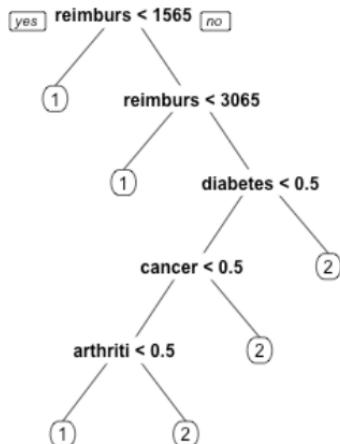
```
%r
```

```
ClaimsTree = rpart(bucket2009 ~ age + alzheimers + arthritis + cancer +
copd + depression + diabetes + heart.failure + ihd + kidney + osteoporosis +
stroke + bucket2008 + reimbursement2008, data=ClaimsTrain,
method="class", cp=0.001)
```

```
%r {"imageWidth": "700px"}
```

```
prp(ClaimsTree)

# The size of the tree depends on the volume of training data and number of
classifications
```



## 12. Make predictions

```
%r
```

```
PredictTest = predict(ClaimsTree, newdata = ClaimsTest, type =
"class")
```

## 13. Confusion Matrix

```
%r
```

```
| confusionMatrixTest = base::table(ClaimsTest$bucket2009,  
| PredictTest)
```

## 14. Prediction Accuracy

```
| %r  
| sum(diag(confusionMatrixTest)) / nrow(ClaimsTest)  
| #Remember baseline accuracy was 0.6838135
```

```
[1] 0.7105545
```

The accuracy of CART model is 71.05%

## 15. Penalty Error

```
| %r  
| as.matrix(base::table(ClaimsTest$bucket2009, PredictTest)) *  
| PenaltyMatrix
```

PredictTest	1	2	3	4	5
1	0	9274	0	0	0
2	36738	0	0	0	0
3	31648	16956	0	0	0
4	18456	19444	0	0	0
5	2752	4278	0	0	0

## 16. Penalty Error Rate on Test Data based on the prediction result

```
| %r  
| sum(as.matrix(base::table(ClaimsTest$bucket2009, PredictTest)) *  
| PenaltyMatrix) / nrow(ClaimsTest)  
| # Remember penalty error rate based on baseline data was 0.7386055
```

```
[1] 0.7617057
```

The penalty error of the baseline method is 0.758

While we increased the accuracy, the penalty error also went up. Why?

By default, rpart will try to maximize the overall accuracy, and every type of error is seen as having a penalty of one.

Our CART model predicts 3, 4, and 5 so rarely because there are very few observations in these classes. Therefore, we do not really expect this model to do better on the penalty error than the baseline method.

To fix this, rpart function allows us to specify a parameter called loss.

### 17. New CART model with loss matrix

```
%r  
ClaimsTree = rpart(bucket2009 ~ age + alzheimers + arthritis +  
cancer + copd + depression + diabetes + heart.failure + ihd + kidney  
+ osteoporosis + stroke + bucket2008 + reimbursement2008,  
data=ClaimsTrain, method="class", cp=0.00005, parms = list(loss =  
PenaltyMatrix))
```

### 18. Redo predictions and penalty error

```
%r  
PredictTest = predict(ClaimsTree, newdata = ClaimsTest, type =  
"class")  
base:::table(ClaimsTest$bucket2009, PredictTest)
```

PredictTest	1	2	3	4	5
1	94310	25295	3087	286	0
2	7176	18942	8079	643	0
3	3590	7706	4692	401	1
4	1304	3193	2803	636	1
5	135	356	408	156	2

### 19. Baseline Accuracy for each bucket

```
%r  
base:::table(ClaimsTest$bucket2009) / nrow(ClaimsTest)
```

0.671269964 0.190172596 0.089464089 0.043323763 0.005769588

### 20. Prediction Accuracy for each bucket

```
%r
```

```

buckets = base::table(ClaimsTest$bucket2009)
accuracy_buckets = diag(base::table(ClaimsTest$bucket2009,
PredictTest)) / buckets
accuracy_buckets

```

0.766885134 0.543685419 0.286272117 0.080131032 0.001892148

## 21. Penalty Error Rate

```

%r
sum(as.matrix(base::table(ClaimsTest$bucket2009, PredictTest)) *
PenaltyMatrix) / nrow(ClaimsTest)

#Remember baseline accuracy was 0.6838135 and accuracy based on
previous model was 0.7578902

```

[1] 0.6418161

The accuracy of this CART model is the 64.73%.

## 22. Penalty Error Rate based on model with penalty condition

```

%r
sum(as.matrix(base::table(ClaimsTest$bucket2009, PredictTest)) *
PenaltyMatrix) / nrow(ClaimsTest)

# Remember penalty error rate based on baseline data was 0.7386055
# And model without penalty condition was 0.7578902

```

[1] 0.6418161

The penalty error of the CART model with loss is 0.642.

Our accuracy is now lower than the baseline method, but our penalty error is also much lower.

### Observations:

Model accuracy of each bucket improved substantially

Penalty Error Rate has gone down

Sacrificed model accuracy to reduce penalty error rate

## Clustering

It is an unsupervised algorithm help segment data into groups

Two popular clustering algorithms –

1. Hierarchical and
2. k-means (works best with large dataset)

Similarity is calculated as distance between two points based on feature values

```
%r #Euclidean distance between 2 movies

God_Father = c(0,1,0,0,0,0,1,0,1,0,0,0,0,0,0,0,0,0,0,0)
Titanic = c(0,1,0,0,0,0,0,1,0,0,0,0,0,1,0,0,0,0,0,0)

#Distance between these two data points in feature space
sqrt(sum((God_Father - Titanic) ^ 2))

#God_Father - Titanic
```

[1] 1.414214

## Hierarchical Clustering

- Start each point with its own cluster.
- Find two nearest clusters based on the Euclidean distance between two centroids and merge into a single cluster.
- Eventually all points will belong to a single cluster.
- The clustering data is visualized using a dendrogram.
- Height presents the distance before the 2 clusters were combined.
- Dendrogram can help us decide how many clustering do we want.

### **Example1:**

1. Load the movie\_lens.txt file and list the variables.

```
%r  
  
movies = read.table("/data/movie_lens.txt", header=FALSE,  
sep="|", quote="\\")  
str(movies)
```

```
'data.frame': 1682 obs. of 24 variables:  
 $ V1 : int 1 2 3 4 5 6 7 8 9 10 ...  
 $ V2 : Factor w/ 1664 levels "'Til There Was You (1997)",...: 1525 618 555 594 344 1318 1545 111 391 1240 ...  
 $ V3 : Factor w/ 241 levels "", "01-Aug-1997",...: 71 71 71 71 71 71 71 71 71 182 ...  
 $ V4 : logi NA NA NA NA NA NA ...  
 $ V5 : Factor w/ 1661 levels "", "<a href='http://us.imdb.com/M/title-exact?%22Langoliers,%20The%22%20(1995)%20(mini)%22,...'>http://us.  
imdb.com/M/title-exact?%22Langoliers,%20The%22%20(1995)%20(mini)",...:</a> 1430 564 504 542 309 1661 1452 102 356 1182 ...  
 $ V6 : int 0 0 0 0 0 0 0 0 0 0 ...  
 $ V7 : int 0 1 0 1 0 0 0 0 0 0 ...  
 $ V8 : int 0 1 0 0 0 0 0 0 0 0 ...  
 $ V9 : int 1 0 0 0 0 0 0 0 0 0 ...  
 $ V10: int 1 0 0 0 0 0 0 1 0 0 ...  
 $ V11: int 1 0 0 1 0 0 0 1 0 0 ...  
 $ V12: int 0 0 0 0 1 0 0 0 0 0 ...  
 $ V13: int 0 0 0 0 0 0 0 0 0 0 ...  
 $ V14: int 0 0 1 1 1 1 1 1 1 1 ...  
 $ V15: int 0 0 0 0 0 0 0 0 0 0 ...
```

2. Add names to the column

```
%r  
  
colnames(movies) = c("ID", "Title", "ReleaseDate",  
"VideoReleaseDate", "IMDB", "Unknown", "Action", "Adventure",  
"Animation", "Childrens", "Comedy", "Crime", "Documentary", "Drama",  
"Fantasy", "FilmNoir", "Horror", "Musical", "Mystery", "Romance",  
"SciFi", "Thriller", "War", "Western")
```

3. List the variables again

```
%r str(movies)
```

```
'data.frame': 1682 obs. of 24 variables:  
$ ID : int 1 2 3 4 5 6 7 8 9 10 ...  
$ Title : Factor w/ 1664 levels "'Til There Was You (1997)",...: 1525 618 555 594 344 1318 1545 111 391 1240 ...  
$ ReleaseDate : Factor w/ 241 levels "", "01-Aug-1997",...: 71 71 71 71 71 71 71 71 71 182 ...  
$ VideoReleaseDate: logi NA NA NA NA NA ...  
$ IMDB : Factor w/ 1661 levels "", <a href="http://us.imdb.com/M/title-exact%22Langoliers,%20The%22%20(1995)%20(mini)%22,...>">http://us.imdb.com/M/title-exact%22Langoliers,%20The%22%20(1995)%20(mini)"..., </a> 1430 564 504 542 309 1661 1452 102 356 1182 ...  
$ Unknown : int 0 0 0 0 0 0 0 0 0 0 ...  
$ Action : int 0 1 0 1 0 0 0 0 0 0 ...  
$ Adventure : int 0 1 0 0 0 0 0 0 0 0 ...  
$ Animation : int 1 0 0 0 0 0 0 0 0 0 ...  
$ Childrens : int 1 0 0 0 0 0 0 1 0 0 ...  
$ Comedy : int 1 0 0 1 0 0 0 1 0 0 ...  
$ Crime : int 0 0 0 0 1 0 0 0 0 0 ...  
$ Documentary : int 0 0 0 0 0 0 0 0 0 0 ...  
$ Drama : int 0 0 0 1 1 1 1 1 1 1 ...  
$ Fantasy : int 0 0 0 0 0 0 0 0 0 0 ...  
$ FilmNoir : int 0 0 0 0 0 0 0 0 0 0 ...  
$ Horror : int 0 0 0 0 0 0 0 0 0 0 ...  
$ Musical : int 0 0 0 0 0 0 0 0 0 0 ...  
$ Mystery : int 0 0 0 0 0 0 0 0 0 0 ...  
$ Romance : int 0 0 0 0 0 0 0 0 0 0 ...  
$ SciFi : int 0 0 0 0 0 0 0 0 0 0 ...
```

#### 4. Remove unnecessary variables

```
%r  
  
movies$ID = NULL  
  
movies$ReleaseDate = NULL  
  
movies$VideoReleaseDate = NULL  
  
movies$IMDB = NULL  
  
  
str(movies)
```

```
'data.frame': 1682 obs. of 20 variables:  
$ Title : Factor w/ 1664 levels "'Til There Was You (1997)",...: 1525 618 555 594 344 1318 1545 111 391 1240 ...  
$ Unknown : int 0 0 0 0 0 0 0 0 0 0 ...  
$ Action : int 0 1 0 1 0 0 0 0 0 0 ...  
$ Adventure : int 0 1 0 0 0 0 0 0 0 0 ...  
$ Animation : int 1 0 0 0 0 0 0 0 0 0 ...  
$ Childrens : int 1 0 0 0 0 0 0 1 0 0 ...  
$ Comedy : int 1 0 0 1 0 0 0 1 0 0 ...  
$ Crime : int 0 0 0 0 1 0 0 0 0 0 ...  
$ Documentary: int 0 0 0 0 0 0 0 0 0 0 ...  
$ Drama : int 0 0 0 1 1 1 1 1 1 1 ...  
$ Fantasy : int 0 0 0 0 0 0 0 0 0 0 ...  
$ FilmNoir : int 0 0 0 0 0 0 0 0 0 0 ...  
$ Horror : int 0 0 0 0 0 0 0 0 0 0 ...  
$ Musical : int 0 0 0 0 0 0 0 0 0 0 ...  
$ Mystery : int 0 0 0 0 0 0 0 0 0 0 ...  
$ Romance : int 0 0 0 0 0 0 0 0 0 0 ...  
$ SciFi : int 0 0 0 0 0 0 0 0 0 0 ...
```

#### 5. Remove duplicate entries from the dataset

```
%r # Remove duplicates  
  
movies = unique(movies)  
  
str(movies)
```

```
'data.frame': 1664 obs. of 20 variables:  
$ Title : Factor w/ 1664 levels "'Til There Was You (1997)",...: 1525 618 555 594 344 1318 1545 111 391 1240 ...  
$ Unknown : int 0 0 0 0 0 0 0 0 0 ...  
$ Action : int 0 1 0 1 0 0 0 0 0 0 ...  
$ Adventure : int 0 1 0 0 0 0 0 0 0 0 ...  
$ Animation : int 1 0 0 0 0 0 0 0 0 ...  
$ Childrens : int 1 0 0 0 0 0 0 1 0 0 ...  
$ Comedy : int 1 0 0 1 0 0 0 1 0 0 ...  
$ Crime : int 0 0 0 0 1 0 0 0 0 0 ...  
$ Documentary: int 0 0 0 0 0 0 0 0 0 ...  
$ Drama : int 0 0 0 1 1 1 1 1 1 1 ...  
$ Fantasy : int 0 0 0 0 0 0 0 0 0 0 ...  
$ FilmNoir : int 0 0 0 0 0 0 0 0 0 0 ...  
$ Horror : int 0 0 0 0 0 0 0 0 0 0 ...  
$ Musical : int 0 0 0 0 0 0 0 0 0 0 ...  
$ Mystery : int 0 0 0 0 0 0 0 0 0 0 ...  
$ Romance : int 0 0 0 0 0 0 0 0 0 0 ...  
$ SciFi : int 0 0 0 0 0 0 0 0 0 0 ...
```

## 6. Compute the euclidean distance

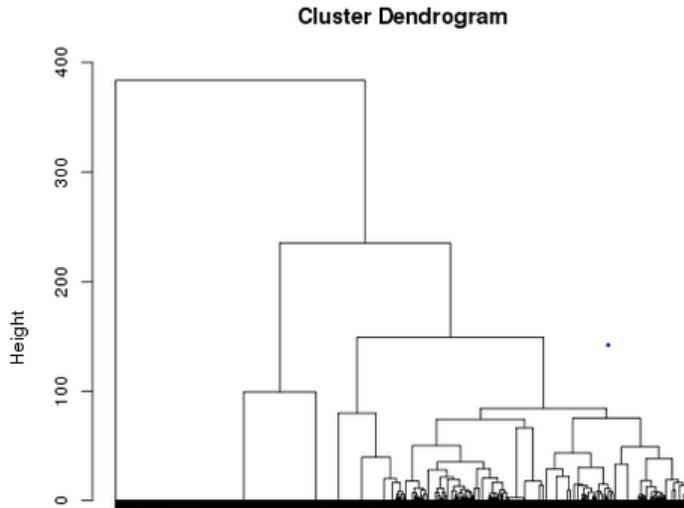
```
%r  
  
distances = dist(movies[2:20], method = "euclidean")  
#Clustering is based on column 2 through 20.
```

## 7. Calculates distance using centroid and variance within a cluster.

```
%r  
  
# Hierarchical clustering  
clusterMovies = hclust(distances, method = "ward.D")  
#ward.D
```

## 8. Plot the Dendrogram

```
%r  
  
plot(clusterMovies)
```



## 9. Assign points to cluster

```
%r
clusterGroups = cutree(clusterMovies, k = 10)
```

```
%r clusterGroups
```

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
1	2	3	1	3	4	1	1	4	1	3	3	5	6	4
16	17	18	19	20	21	22	23	24	25	26	27	28	29	30
7	3	4	4	6	1	1	3	3	5	5	2	2	3	4
31	32	33	34	35	36	37	38	39	40	41	42	43	44	45
1	8	6	9	1	6	4	2	3	5	5	5	3	3	9
46	47	48	49	50	51	52	53	54	55	56	57	58	59	60
4	9	8	7	1	1	4	2	2	6	3	4	4	4	4
61	62	63	64	65	66	67	68	69	70	71	72	73	74	75
4	2	1	4	9	7	5	6	1	7	1	1	1	1	8
76	77	78	79	80	81	82	83	84	85	86	87	88	89	90
3	3	1	2	1	7	2	7	2	5	4	4	7	3	6
91	92	93	94	95	96	97	98	99	100	101	102	103	104	105
1	6	9	1	1	2	1	3	1	3	2	1	1	5	5
106	107	108	109	110	111	112	113	114	115	116	117	118	119	120
3	4	5	2	1	7	1	4	1	8	5	2	2	8	1
121	122	123	124	125	126	127	128	129	130	131	132	133	134	135

## 10. Let's use the tapply function to compute the percentage of movies in each genre and cluster

```
%r
# Now let's figure out what the clusters are like.
```

```
tapply(movies$Action, clusterGroups, mean) #This will calculate mean  
of Action variable for each of 10 clusters
```

```
0.1784512 0.7839196 0.1238532 0.0000000 0.0000000 0.1015625 0.0000000  
8 9 10  
0.0000000 0.0000000 0.0000000
```

```
%r tapply(movies$Romance, clusterGroups, mean)
```

```
0.10437710 0.04522613 0.03669725 0.00000000 0.00000000 1.00000000  
7 8 9 10  
1.00000000 0.00000000 0.00000000 0.00000000
```

```
%r  
  
clusters = data.frame()  
  
for(i in 1:10){  
  
    clusters = rbind(clusters, colMeans(subset(movies[2:20],  
clusterGroups == i)))  
  
}  
  
clusters = t(clusters)  
  
rownames(clusters) = colnames(movies[2:20])  
  
clusters
```

Genre	Mean 1	Mean 2	Mean 3	Mean 4	Mean 5	Mean 6	Mean 7	Mean 8	Mean 9	Mean 10
Unknown	0.006734007	0.000000000	0.000000000	0	0	0.0000000	0	0.0000000	0	0.0000000
Action	0.178451178	0.783919598	0.123853211	0	0	0.1015625	0	0.0000000	0	0.0000000
Adventure	0.185185185	0.351758794	0.036697248	0	0	0.0000000	0	0.0000000	0	0.0000000
Animation	0.134680135	0.010050251	0.000000000	0	0	0.0000000	0	0.0000000	0	0.0000000
Childrens	0.393939394	0.005025126	0.009174312	0	0	0.0000000	0	0.0000000	0	0.0000000
Comedy	0.363636364	0.065326633	0.064220183	0	1	0.1093750	1	0.0212766	0	0.0000000
Crime	0.033670034	0.005025126	0.412844037	0	0	0.0468750	0	0.0000000	0	0.0000000
Documentary	0.010101010	0.000000000	0.000000000	0	0	0.0000000	0	0.0000000	0	1.0000000
Drama	0.306397306	0.110552764	0.380733945	1	0	0.6640625	0	0.0000000	0	0.0000000
Fantasy	0.070707071	0.000000000	0.004587156	0	0	0.0000000	0	0.0000000	0	0.0000000
FilmNoir	0.000000000	0.000000000	0.105504587	0	0	0.0078125	0	0.0000000	0	0.0000000
Horror	0.016835017	0.080402010	0.018348624	0	0	0.0156250	0	0.0000000	0	0.0000000
Musical	0.188552189	0.000000000	0	0	0.0000000	0	0.0000000	0	0.0000000	0
Mystery	0.000000000	0.000000000	0.275229358	0	0	0.0000000	0	0.0000000	0	0.0000000
Romance	0.104377104	0.045226131	0.036697248	0	0	1.0000000	1	0.0000000	0	0.0000000
SciFi	0.074074074	0.346733668	0.041284404	0	0	0.0000000	0	0.0000000	0	0.0000000
Thriller	0.040404040	0.376884422	0.610091743	0	0	0.1406250	0	0.0000000	0	0.0000000

11. Find which cluster 'Men in Black' is in.

```
%r  
  
subset(movies, Title=="Men in Black (1997)")
```

```

257 Men in Black (1997)      0      1      1      0      0
Comedy Crime Documentary Drama Fantasy FilmNoir Horror Musical Mystery
257      1      0      0      0      0      0      0      0
Romance SciFi Thriller War Western
257      0      1      0      0      0

```

```

%r
clusterGroups[which(movies$title=="Men in Black (1997)")]

```

```

257
2

```

12. Create a new data set with just the movies from cluster 2

```

%r
cluster2 = subset(movies, clusterGroups==2)

```

13. Look at the first 10 titles in this cluster

```

%r
cluster2$title[1:10]

```

```

[1] GoldenEye (1995)<br />
[2] Bad Boys (1995)<br />
[3] Apollo 13 (1995)<br />
[4] Net, The (1995)<br />
[5] Natural Born Killers (1994)<br />
[6] Outbreak (1995)<br />
[7] Stargate (1994)<br />
[8] Fugitive, The (1993)<br />
[9] Jurassic Park (1993)<br />
[10] Robert A. Heinlein's The Puppet Masters (1994)
1664 Levels: 'Til There Was You (1997) ...

```

### **Example2:**

1. Load the flower.csv file and list the variables

```

%r
flower = read.csv("/data/flower.csv", header=FALSE)
str(flower)

```

```
'data.frame': 50 obs. of 50 variables:  
 $ V1 : num 0.0991 0.0991 0.1034 0.1034 0.1034 ...  
 $ V2 : num 0.112 0.108 0.112 0.116 0.108 ...  
 $ V3 : num 0.134 0.116 0.121 0.116 0.112 ...  
 $ V4 : num 0.138 0.138 0.121 0.121 0.112 ...  
 $ V5 : num 0.138 0.134 0.125 0.116 0.112 ...  
 $ V6 : num 0.138 0.129 0.121 0.108 0.112 ...  
 $ V7 : num 0.129 0.116 0.103 0.108 0.112 ...  
 $ V8 : num 0.116 0.103 0.103 0.103 0.116 ...  
 $ V9 : num 0.1121 0.0991 0.1078 0.1121 0.1164 ...  
 $ V10: num 0.121 0.108 0.112 0.116 0.125 ...  
 $ V11: num 0.134 0.125 0.129 0.134 0.129 ...  
 $ V12: num 0.147 0.134 0.138 0.129 0.138 ...  
 $ V13: num 0.000862 0.146552 0.142241 0.142241 0.133621 ...  
 $ V14: num 0.000862 0.000862 0.142241 0.133621 0.12931 ...  
 $ V15: num 0.142 0.142 0.134 0.121 0.116 ...  
 $ V16: num 0.125 0.125 0.116 0.108 0.108 ...
```

## 2. Change the data type to matrix

```
%r  
  
flowerMatrix = as.matrix(flower)  
  
str(flowerMatrix)
```

```
num [1:50, 1:50] 0.0991 0.0991 0.1034 0.1034 0.1034 ...
```

## 3. Turn matrix into a vector

```
%r  
  
flowerVector = as.vector(flowerMatrix)  
  
str(flowerVector)
```

```
num [1:2500] 0.0991 0.0991 0.1034 0.1034 0.1034 ...
```

## 4. Compute the euclidean distance

```
%r  
  
distance = dist(flowerVector, method = "euclidean")
```

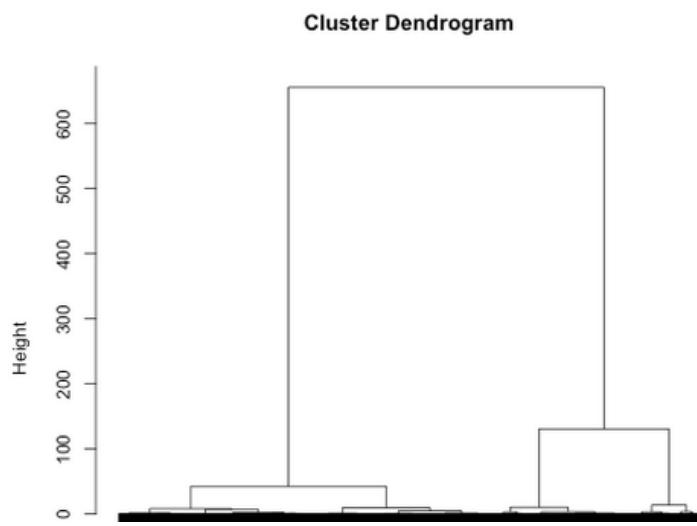
## 5. Prepare the data for clustering

```
%r  
  
# Hierarchical clustering  
clusterIntensity = hclust(distance, method="ward.D")
```

## 6. Plot the dendrogram

```
%r
```

```
| plot(clusterIntensity)
```

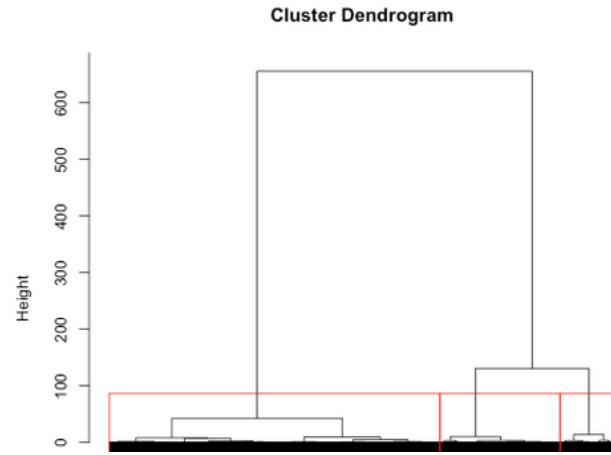


7. Select 3 cluster and plot the dendrogram

```
| %r
```

```
| plot(clusterIntensity)
```

```
| rect.hclust(clusterIntensity, k = 3, border = "red")
```



```
| %r
```

```
| flowerClusters = cutree(clusterIntensity, k = 3)
```

```
| flowerClusters
```

```
[1] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
[35] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
[69] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
[103] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
[137] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
[171] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
[205] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
[239] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
[273] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
[307] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
[341] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 3 3  
[375] 2 2 3 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
[409] 1 1 1 1 1 1 1 1 1 2 2 1 1 1 3 3 3 3 3 2 1 2 3 2 1 1 1 1 1 1 1 1  
[443] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 3 3 1 1 1 3 3 3  
[477] 3 2 2 3 3 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
[511] 1 1 1 1 1 1 2 3 3 2 1 1 3 3 3 3 2 3 3 3 1 2 3 3 3 1 1 1 1 1 1 1  
[545] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 3 2 1 2 3 3 3 1 1 3 3 3 3 2
```

#### 8. Find the mean intensity values using tapply function

```
%r  
tapply(flowerVector, flowerClusters, mean)
```

```
0.08574315 0.50826255 0.93147713
```

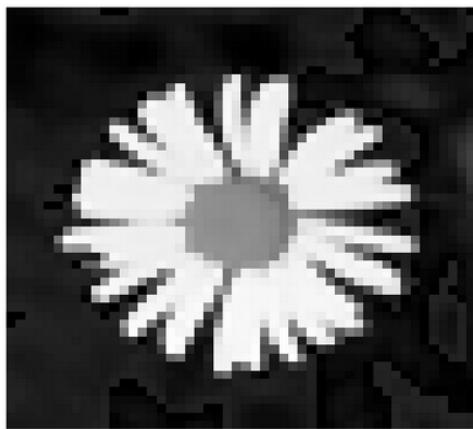
#### 9. Plot the image and the clusters

```
%r  
dim(flowerClusters) = c(50,50)  
image(flowerClusters, axes = FALSE)
```



```
%r  
# Original image
```

```
| image(flowerMatrix, axes=FALSE, col=grey(seq(0,1,length=256)))
```



## k-means Clustering

\*Algorithm in nutshell\*

Step 1: Specify a desired number of clusters k based on previous knowledge or experimenting

Step 2: Randomly assign each data to a cluster

Step 3: Compute cluster centroids

Step 4: Reassign each point to the closest cluster centroid

Step 5: Repeat 4 and 5 until no improvement is made

## MRI Image Clustering using K-means

1. Load the tumor.csv file and list the variables

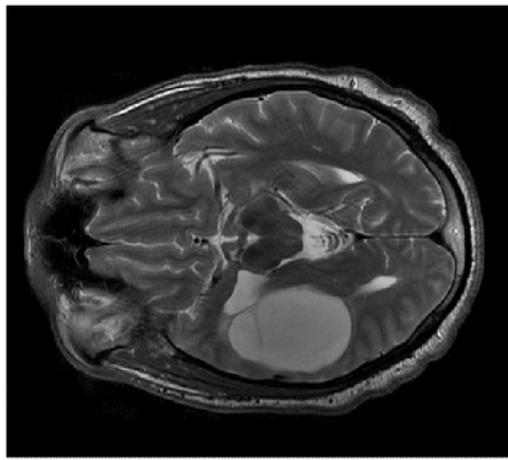
```
%r  
healthy = read.csv("/data/tumor.csv", header=FALSE)  
healthyMatrix = as.matrix(healthy)  
str(healthyMatrix)
```

num [1:571, 1:512] 0 0 0 0 0 0 0 0 ...

2. Plot the image

```
%r
```

```
| image(healthyMatrix,axes=FALSE,col=grey(seq(0,1,length=256)))
```



### 3. Run k-means algorithm with number of clusters as 5

```
%r  
# Run k-means  
set.seed(1)  
healthyVector = as.vector(healthyMatrix)  
KMC = kmeans(healthyVector, centers = k, iter.max = 1000)  
str(KMC)
```

```
List of 9  
$ cluster : int [1:292352] 5 5 5 5 5 5 5 5 ...  
$ centers : num [1:5, 1] 0.43396 0.61961 0.28868 0.18492 0.00708  
... attr(*, "dimnames")=List of 2  
... ..$ : chr [1:5] "1" "2" "3" "4" ...  
... ..$ : NULL  
$ totss : num 8600  
$ withinss : num [1:5] 54.3 50.6 69 76.8 43.8  
$ tot.withinss: num 294  
$ betweenss : num 8305  
$ size : int [1:5] 29811 8366 60806 53943 139426  
$ iter : int 5  
$ ifault : int 0
```

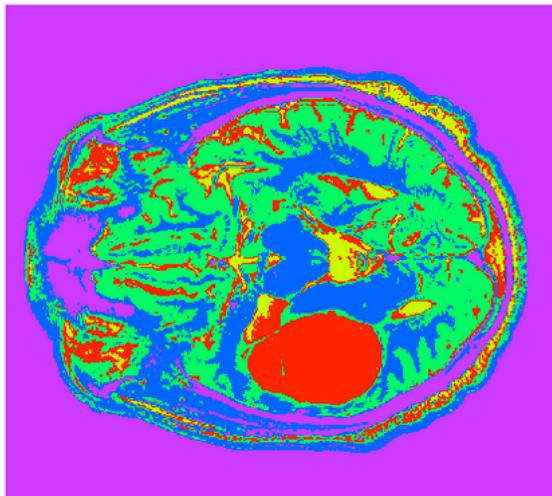
### 4. Extract clusters

```
%r  
healthyClusters = KMC$cluster  
KMC$centers[2]
```

```
[1] 0.6196116
```

## 5. Plot the image of clusters

```
%r  
dim(healthyClusters) = c(nrow(healthyMatrix), ncol(healthyMatrix))  
image(healthyClusters, axes = FALSE, col=rainbow(k))
```



## PCA (Principle Component Analysis)

```
%r  
library(ISLR)  
nci.labs=NCI60$labs  
nci.data=NCI60$data  
dim(nci.data)
```

[1] 64 6830

```
%r pr.out=prcomp(nci.data, scale=TRUE) #Performs a principal  
components analysis on the given data matrix and returns the results  
as an object of class prcomp.  
pr.out
```

```
Standard deviations:  
[1] 2.785347e+01 2.148136e+01 1.982046e+01 1.703256e+01 1.597181e+01  
[6] 1.572108e+01 1.447145e+01 1.354427e+01 1.314400e+01 1.273860e+01  
[11] 1.268672e+01 1.215769e+01 1.183019e+01 1.162554e+01 1.143779e+01  
[16] 1.100051e+01 1.065666e+01 1.048880e+01 1.043518e+01 1.032194e+01  
[21] 1.014608e+01 1.005439e+01 9.982655e+00 9.647656e+00 9.507638e+00  
[26] 9.332529e+00 9.273200e+00 9.090046e+00 8.981173e+00 8.750031e+00  
[31] 8.599622e+00 8.447375e+00 8.373048e+00 8.215787e+00 8.157313e+00  
[36] 7.974655e+00 7.904462e+00 7.821271e+00 7.721562e+00 7.586035e+00  
[41] 7.456193e+00 7.344380e+00 7.104489e+00 7.013055e+00 6.958385e+00  
[46] 6.866265e+00 6.807439e+00 6.647630e+00 6.616068e+00 6.407926e+00  
[51] 6.219838e+00 6.203258e+00 6.067065e+00 5.918049e+00 5.912333e+00  
[56] 5.735386e+00 5.472610e+00 5.292148e+00 5.021174e+00 4.683979e+00  
[61] 4.175673e+00 4.082121e+00 4.041243e+00 2.148247e-14
```

```
%r str(pr.out$x)
```

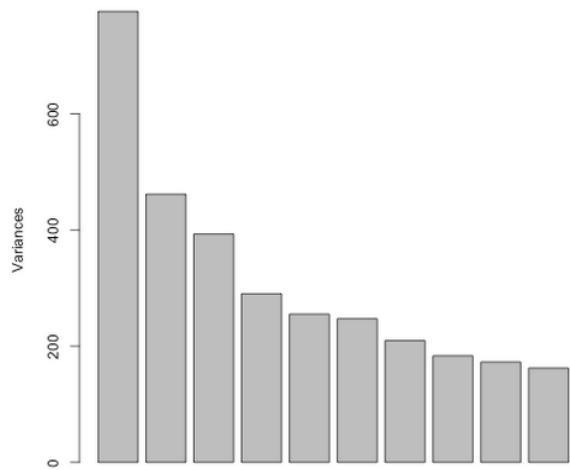
```
num [1:64, 1:64] -19.7 -22.9 -27.2 -42.5 -55 ...
```

```
%r capture.output(summary(pr.out))
```

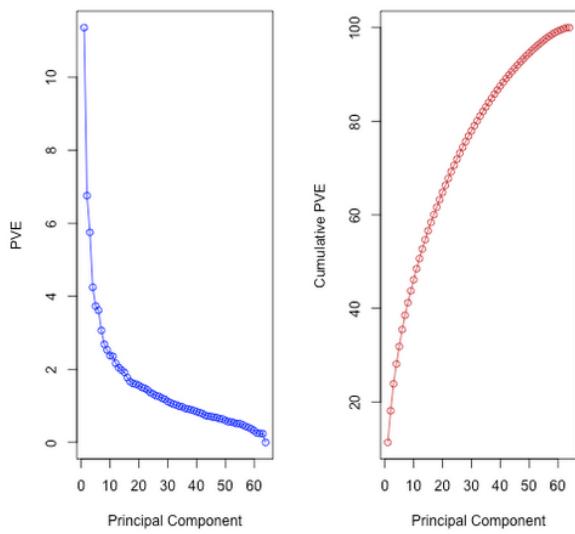
```
[1] "Importance of components:<br />"  
[2] " PC1      PC2      PC3      PC4      PC5"<br />  
[3] "Standard deviation    27.8535 21.48136 19.82046 17.03256 15.97181"<br />  
[4] "Proportion of Variance 0.1136 0.06756 0.05752 0.04248 0.03735"<br />  
[5] "Cumulative Proportion 0.1136 0.18115 0.23867 0.28115 0.31850"<br />  
[6] " PC6      PC7      PC8      PC9      PC10"<br />  
[7] "Standard deviation    15.72108 14.47145 13.54427 13.14400 12.73860"<br />  
[8] "Proportion of Variance 0.03619 0.03066 0.02686 0.02529 0.02376"<br />  
[9] "Cumulative Proportion 0.35468 0.38534 0.41220 0.43750 0.46126"<br />  
[10] " PC11     PC12     PC13     PC14     PC15"<br />  
[11] "Standard deviation    12.68672 12.15769 11.83019 11.62554 11.43779"<br />  
[12] "Proportion of Variance 0.02357 0.02164 0.02049 0.01979 0.01915"<br />  
[13] "Cumulative Proportion 0.48482 0.50646 0.52695 0.54674 0.56590"<br />  
[14] " PC16     PC17     PC18     PC19     PC20"<br />  
[15] "Standard deviation    11.00051 10.65666 10.48880 10.43518 10.3219"<br />  
[16] "Proportion of Variance 0.01772 0.01663 0.01611 0.01594 0.0156"<br />  
[17] "Cumulative Proportion 0.58361 0.60024 0.61635 0.63229 0.6479"<br />  
[18] " PC21     PC22     PC23     PC24     PC25     PC26"
```

```
%r plot(pr.out)
```

```
pr.out
```



```
%r  
pve=100*pr.out$sdev^2/sum(pr.out$sdev^2)  
par(mfrow=c(1,2))  
plot(pve, type="o", ylab="PVE", xlab="Principal Component",  
col="blue")  
plot(cumsum(pve), type="o", ylab="Cumulative PVE", xlab="Principal  
Component", col="brown3")
```



```
%r base::table(nci.labs)
```

```

nci.labs
  BREAST      CNS      COLON K562A-repro K562B-repro LEUKEMIA
    7          5         7          1          1          6
MCF7A-repro MCF7D-repro MELANOMA   NSCLC     OVARIAN   PROSTATE
    1          1         8         9         6          2
  RENAL      UNKNOWN
    9          1

```

```

|r
|Cols=function(vec) {
|  cols=rainbow(length(unique(vec))) #rainbow: Create a vector of n
|  contiguous colors.
|
|  return(cols[as.numeric(as.factor(vec))])
|
|}

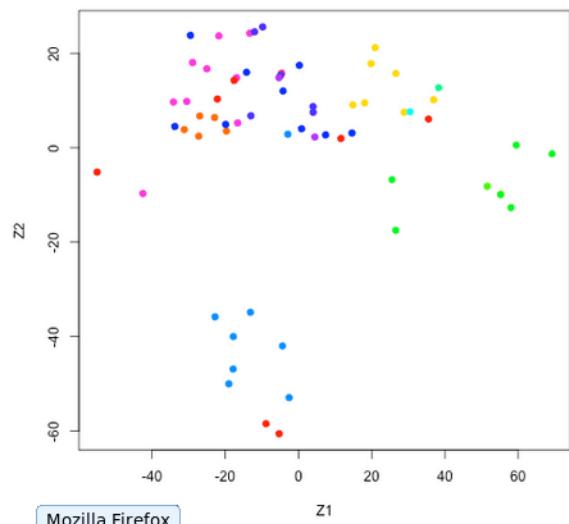
par(mfrow=c(1,2))

```

```

|r plot(pr.out$x[,1:2], col=Cols(nci.labs),
|  pch=19, xlab="Z1", ylab="Z2")

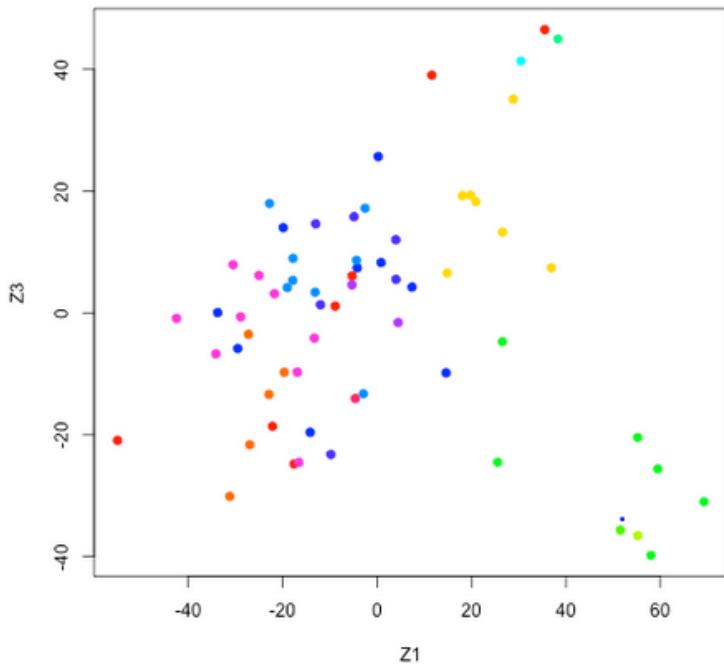
```



```

|r plot(pr.out$x[,c(1,3)], col=Cols(nci.labs),
|  pch=19, xlab="Z1", ylab="Z3")

```



## Exercise10: Graphx

### Analyzing Flight Data

Input: <http://stat-computing.org/dataexpo/2009/the-data.html>

1. Load the data

```

import org.apache.spark.graphx._
import org.apache.spark.rdd.RDD
import scala.util.MurmurHash
val df_1 =
  sqlContext.read.format("com.databricks.spark.csv").option("header",
  "true").load("/user/technocracy/data/2008.csv")

```

```

import org.apache.spark.graphx._
import org.apache.spark.rdd.RDD
import scala.util.MurmurHash
df_1: org.apache.spark.sql.DataFrame = [Year: string, Month: string, DayofMonth: string, DayOfWeek: string, DepTime: string, CRSDepTime: string, ArrTime: string, CRSArrTime: string, UniqueCarrier: string, FlightNum: string, TailNum: string, ActualElapsedTime: string, CRSElapsedTime: string, AirTime: string, ArrDelay: string, DepDelay: string, Origin: string, Dest: string, Distance: string, TaxiIn: string, TaxiOut: string, Cancelled: string, CancellationCode: string, Diverted: string, CarrierDelay: string, WeatherDelay: string, NASDelay: string, SecurityDelay: string, LateAircraftDelay: string]

```

## 2. Check the data

```
df_1.take(2)
```

```
res17: Array[org.apache.spark.sql.Row] = Array([2008,1,3,4,2003,1955,2211,2225,WN,335,N712SW,128,150,116,-14,8,IAD,TPA,810,4,8,0,,0,NA,NA,NA,NA,NA], [2008,1,3,4,754,735,1002,1000,WN,3231,N772SW,128,145,113,2,19,IAD,TPA,810,5,10,0,,0,NA,NA,NA,NA,NA])
```

## 3. Preparing Edge and vertices

```

val flightsFromTo = df_1.select($"Origin", $"Dest")
val airportCodes = df_1.select($"Origin", $"Dest").flatMap(x =>
  Iterable(x(0).toString, x(1).toString))

```

```

flightsFromTo: org.apache.spark.sql.DataFrame = [Origin: string, Dest: string]
airportCodes: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[15] at flatMap at <console>:44

```

- Convert our airport codes to primitive data types, make them distinct, and turn them into a list of vertices. That involves two conceptual maps. The distinct values of what we have now, mapped to a unique integer id, then turned into a list of vertices.

```

val airportVertices: RDD[(VertexId, String)] =
  airportCodes.distinct().map(x => (MurmurHash.stringHash(x), x))
val defaultAirport = ("Missing")

```

```

warning: there were 1 deprecation warning(s); re-run with -deprecation for details
airportVertices: org.apache.spark.rdd.RDD[(org.apache.spark.graphx.VertexId, String)] = MapPartitionsRDD[19] at map at <console>:46
defaultAirport: String = Missing

```

- Count the number of flights to and from airports, convert them into edges from airport ID to airport ID (of the primitive type), then validate and convert them with the Edgecase class.

```
val flightEdges = flightsFromTo.map(x =>
  (MurmurHash.stringHash(x(0).toString), MurmurHash.stringHash(x(1).toString)), 1)).reduceByKey(_ + _).map(x => Edge(x._1._1, x._1._2, x._2))
```

```
warning: there were 2 deprecation warning(s); re-run with -deprecation for details
flightEdges: org.apache.spark.rdd.RDD[org.apache.spark.graphx.Edge[Int]] = MapPartitionsRDD[28] at map at <console>:47
```

## 6. Create our Graph

```
val graph = Graph(airportVertices, flightEdges, defaultAirport)
graph.persist() // we're going to be using it a lot
```

```
graph: org.apache.spark.graphx.Graph[String,Int] = org.apache.spark.graphx.impl.GraphImpl@7d90ba4
res4: org.apache.spark.graphx.Graph[String,Int] = org.apache.spark.graphx.impl.GraphImpl@7d90ba4
```

## 7. To answer basic questions:

- How many airports are there?
- How many unique flights from airport A to airport B are there?

```
graph.numVertices
graph.numEdges
```

```
res6: Long = 305
res7: Long = 5366
```

## 8. What are the top 10 flights from airport to airport?

```
graph.triplets.sortBy(_.attr, ascending=false).map(triplet => "There
were " + triplet.attr.toString + " flights from " + triplet.srcAttr
+ " to " + triplet.dstAttr + ".").take(10)
```

```
res9: Array[String] = Array(There were 13788 flights from SFO to LAX., There were 13390 flights from LAX to SFO., There were 12383 flights from OGG to HNL., There were 120
35 flights from LGA to BOS., There were 12029 flights from BOS to LGA., There were 12014 flights from HNL to OGG., There were 11773 flights from LAX to LAS., There were 11
729 flights from LAS to LAX., There were 11257 flights from LAX to SAN., There were 11224 flights from SAN to LAX.)
```

## 9. What are the lowest 10 flights from airport to airport?

```
|-----|  
| graph.triplets.sortBy(_.attr).map(triplet => "There were " +  
| triplet.attr.toString + " flights from " + triplet.srcAttr + " to "  
| + triplet.dstAttr + ".") .take(10)  
|-----|
```

```
res11: Array[String] = Array(There were 1 flights from RNO to PIH., There were 1 flights from BFL to FAT., There were 1 flights from PDX to MSN., There were 1 flights from DSM to RFD., There were 1 flights from BOI to GJT., There were 1 flights from TVC to MSN., There were 1 flights from LGA to PHL., There were 1 flights from COS to CLE., T  
here were 1 flights from ORF to SAV., There were 1 flights from DEN to ROC.)  
.
```

10. What airport has the most in degrees (which airports are flown to the most and flown out of the most) or unique flights into it?

```
|-----|  
| graph.inDegrees.join(airportVertices).sortBy(_.2._1,  
| ascending=false) .take(1)  
|-----|
```

```
| graph.inDegrees.join(airportVertices).sortBy(_.2._1, ascending=false).take(1)  
res13: Array[(org.apache.spark.graphx.VertexId, (Int, String))] = Array((2042033420,(173,ATL)))
```

and out degree

```
|-----|  
| graph.outDegrees.join(airportVertices).sortBy(_.2._1,  
| ascending=false) .take(1)  
|-----|
```

```
| graph.outDegrees.join(airportVertices).sortBy(_.2._1, ascending=false).take(1)  
res15: Array[(org.apache.spark.graphx.VertexId, (Int, String))] = Array((2042033420,(173,ATL)))
```