

# Hands On: Building a Graph

In this activity, we are going to build our first graph using Spark's GraphX library. You can consult the datasets to be used within the `big-data-5/graphx/EOADATA` directory. The data has already been placed inside the Docker containers that we are going to use.

## Instructions

**Step 1. Set up.** We will start by setting up *Graphx*. First, make sure to open Docker Desktop to have the Docker engine running. Then, open your local terminal shell, and run the following command to download the Docker image for this module.

```
1 docker pull pramonettivega/graphx-course4
```

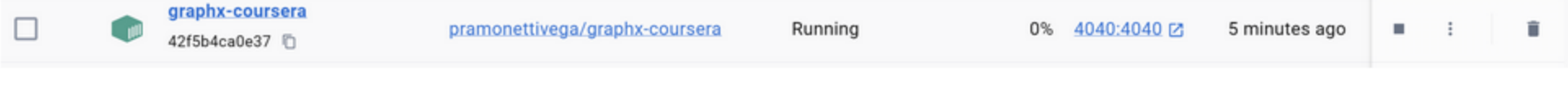
After pulling the image, build your docker container by running the following command:

```
1 docker run -it --name graphx-course4 -p 4040:4040 pramonettivega/graphx-course4
```

After running command, you should be immediately redirected to the container's shell.

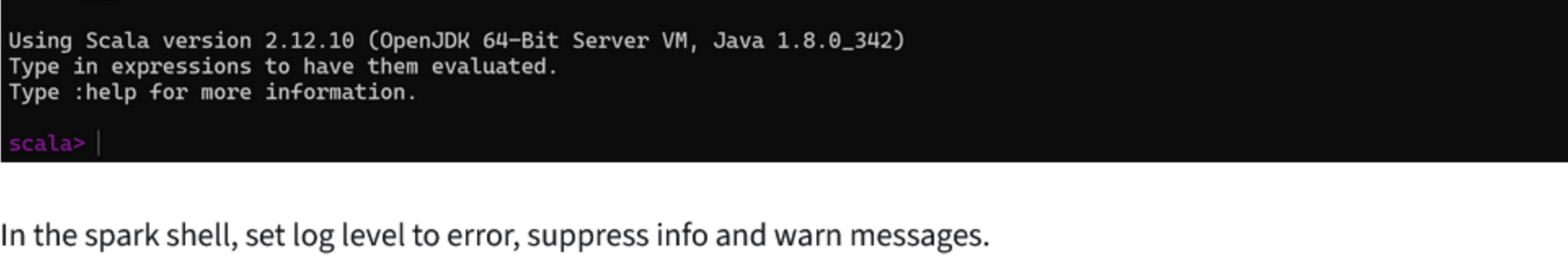


You can also confirm that the container is running.



In the container's shell, start the Spark's shell:

```
1 spark-shell --jars lib/gs-core-1.2.jar,lib/gs-ui-1.2.jar,lib/jcommon1.0.16.jar,lib/jfreechart-1.0.13.jar
```



In the spark shell, set log level to error, suppress info and warn messages.

```
1 import org.apache.log4j.Logger
2 import org.apache.log4j.Level
3 Logger.getLogger("org").setLevel(Level.ERROR)
4 Logger.getLogger("akka").setLevel(Level.ERROR)
5
```

Import the Spark's GraphX and RDD libraries along with Scala's source library.

```
1 import org.apache.spark.graphx._
2 import org.apache.spark.rdd._
3 import scala.io.Source
```

## Step 2. Import the vertices.

Before importing any datasets, let view what the files contain. Print the first 5 lines of each CSV file.

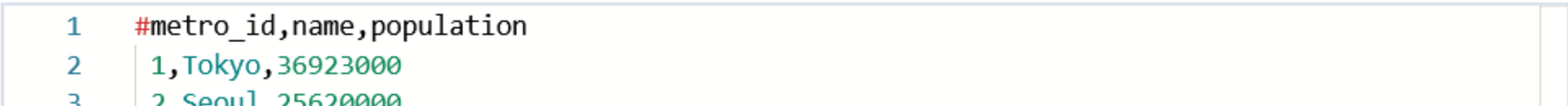
Input:

```
1 Source.fromFile("./EOADATA/metro.csv").getLines().take(5).foreach(println)
```



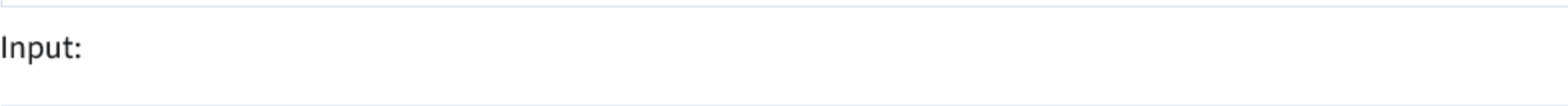
Output:

```
1 #metro_id,name,population
2 1,Tokyo,36923800
3 2,Seoul,25620000
4 3,Shanghai,24750000
5 4,Guangzhou,23900000
6
```



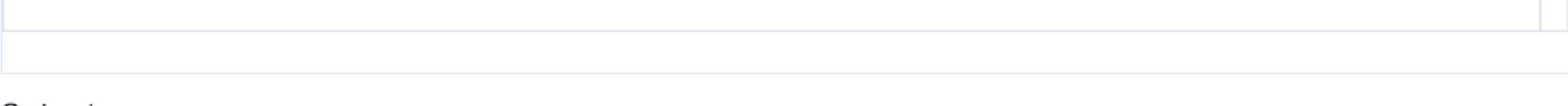
Input:

```
1 Source.fromFile("./EOADATA/country.csv").getLines().take(5).foreach(println)
```



Output:

```
1 #country_id,name
2 1,Japan
3 2,South Korea
4 3,China
5 4,India
6
```



Input:

```
1 Source.fromFile("./EOADATA/metro_country.csv").getLines().take(5).foreach(println)
```



Output:

```
1 #metro_id,country_id
2 1,1
3 2,2
4 3,3
5 4,3
6
```



Now, create case classes for the places (metros and countries).

Input:

```
1 class PlaceNode(val name: String) extends Serializable
```



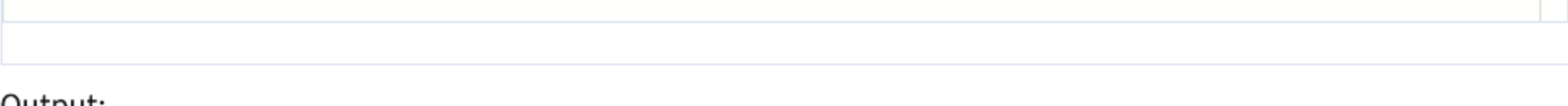
Output:

```
1 defined class PlaceNode
```



Input:

```
1 case class Metro(override val name: String, population: Int) extends PlaceNode(name)
```



Output:

```
1 defined class Metro
```



Input:

```
1 case class Country(override val name: String) extends PlaceNode(name)
```



Output:

```
1 defined class Country
```



Read the CSV file *metros.csv* into an RDD of Metro vertices, ignore lines that start with # and map the columns to: id, Metro(name, population).

Input:

```
1 val metros: RDD[(VertexId, PlaceNode)] =
2   sc.textFile("./EOADATA/metro.csv").
3   filter(l => !l.startsWith("#")).
4   map {line =>
5     val row = line.split(',')
6     (0L + row(0).toInt, Metro(row(1), row(2).toInt))
7   }
```



Output:

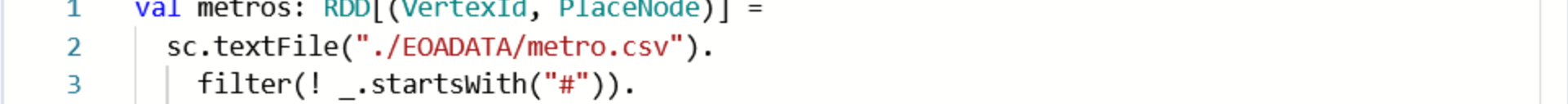
```
1 metros: org.apache.spark.rdd.RDD[(org.apache.spark.graphx.VertexId, PlaceNode)] = MapPartit
```



Read the CSV file *country.csv* into an RDD of Country vertices, ignore lines that start with # and map the columns to: id, Country(name). Add 100 to the country indexes so they are unique from the metro indexes.

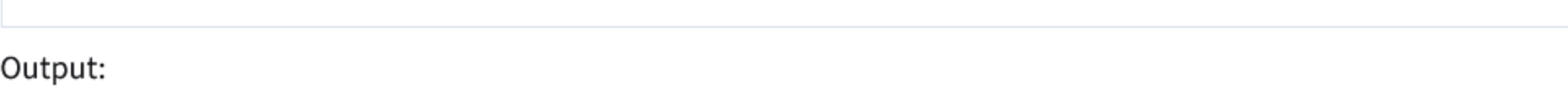
Input:

```
1 val countries: RDD[(VertexId, PlaceNode)] =
2   sc.textFile("./EOADATA/country.csv").
3   filter(l => !l.startsWith("#")).
4   map {line =>
5     val row = line.split(',')
6     (100L + row(0).toInt, Country(row(1)))
7   }
```



Output:

```
1 countries: org.apache.spark.rdd.RDD[(org.apache.spark.graphx.VertexId, PlaceNode)] = MapPartit
```

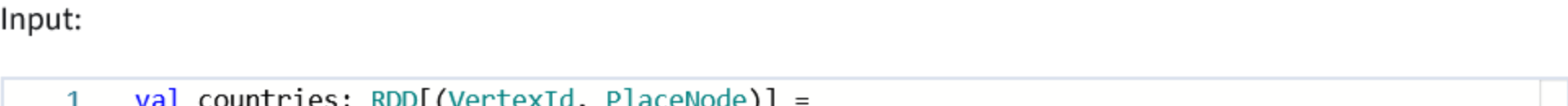


## Step 3. Import the edges.

Read the comma delimited text file *metro\_country.tsv* into an RDD[Edge[Int]] collection. Remember to add 100 to the countries' vertex id.

Input:

```
1 val mclinks: RDD[Edge[Int]] =
2   sc.textFile("./EOADATA/metro_country.csv").
3   filter(l => !l.startsWith("#")).
4   map {line =>
5     val row = line.split(',')
6     Edge(0L + row(0).toInt, 100L + row(1).toInt, 1)
7   }
```



Output:

```
1 mclinks: org.apache.spark.rdd.RDD[org.apache.spark.graphx.Edge[Int]] = MapPartitionsRDD[11]
```



## Step 4. Create a graph.

Concatenate the two sets of nodes into a single RDD.

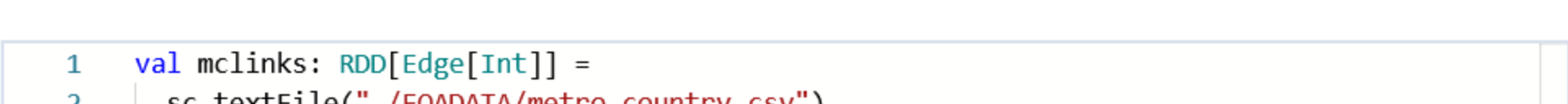
Input:

```
1 val nodes = metros ++ countries
```



Output:

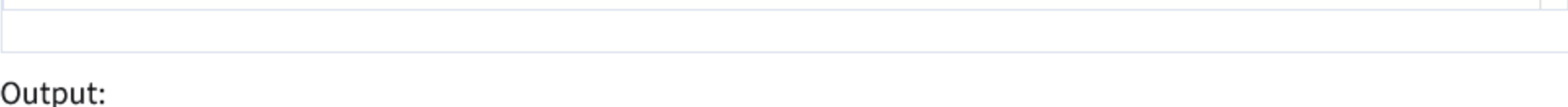
```
1 nodes: org.apache.spark.rdd.RDD[(org.apache.spark.graphx.VertexId, PlaceNode)] = UnionRDD[12]
```



Pass the concatenated RDD to the Graph() factory method along with the RDD link

Input:

```
1 val metrosGraph = Graph(nodes, mclinks)
```



Output:

```
1 metrosGraph: org.apache.spark.graphx.Graph[PlaceNode,Int] = org.apache.spark.graphx.impl.Gra
```



Print the first 5 vertices and edges.

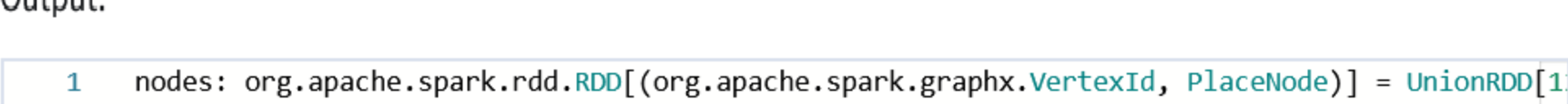
Input:

```
1 metrosGraph.vertices.take(5)
```



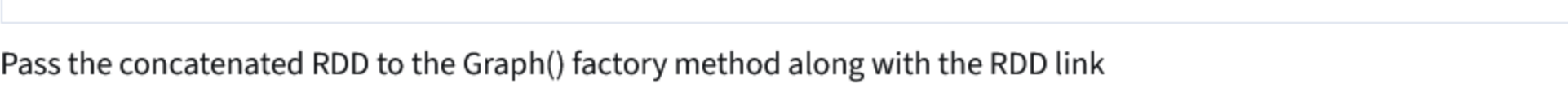
Output:

```
1 res5: Array[(org.apache.spark.graphx.VertexId, PlaceNode)] = Array((52,Metro(Ankara,5150072
2 (56,Metro(Boston,4732161)), (4,Metro(Guangzhou,23900000)), (112,Country(United Kingdom)),
3 (120,Country(Colombia)))
```



Input:

```
1 metrosGraph.edges.take(5)
```



Output:

```
1 res6: Array[org.apache.spark.graphx.Edge[Int]] = Array(Edge(1,101,1), Edge(2,102,1),
2 Edge(3,103,1), Edge(4,103,1), Edge(5,104,1))
```

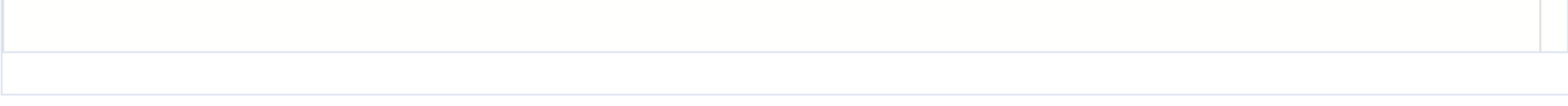


## Step 5. Using Spark's filter method to return Vertices in the graph.

Filter all of the edges in *metrosGraph* that have a source vertex id of 1 and create a map of destination vertex ids.

Input:

```
1 metrosGraph.edges.filter(_._srcId == 1).map(_._dstId).collect()
```



Output:

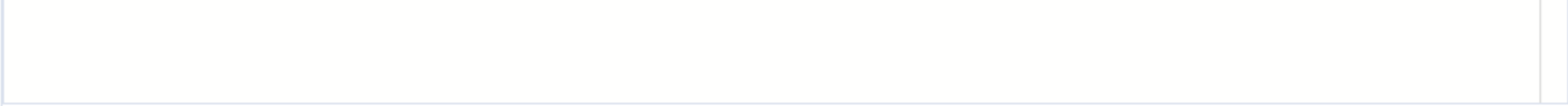
```
1 res7: Array[org.apache.spark.graphx.VertexId] = Array(101)
```



Similarly, filter all of the edges in *metrosGraph* where the destination vertexid is 103 and create a map of all of the source ids.

Input:

```
1 metrosGraph.edges.filter(_._dstId == 103).map(_._srcId).collect()
```



Output:

```
1 res11: Array[org.apache.spark.graphx.VertexId] = Array(3, 4, 7, 24, 34)
```



We recommend you to go directly to the next activity *Hands-On: Building a Degree Histogram*, because that activity will use the graph you created in this activity.

Mark as completed