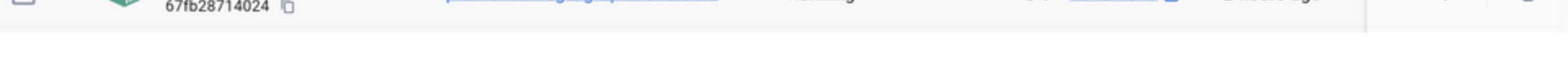


# Hands On: Joining Graph Datasets

In this activity, we are going to:

1. Create a new dataset
2. Join two datasets with `JoinVertices`
3. Join two datasets with `outerJoinVertices`
4. Create a new return type for the joined vertices

**Step 1. Set up.** First, make sure your Docker container is running:



Once you verify the Docker container is running, run the following command to access the container's shell:

```
1 docker exec -it graphx-coursera /bin/sh
```

Once the container's shell is open, start the Spark Shell.

```
1 spark-shell
```

It may take several seconds for the Spark Shell to start. Be patient and wait for the **scala>** prompt.

Set log level to error in order to suppress the info and warn messages so the output is easier to read.

```
1 import org.apache.log4j.Logger
2 import org.apache.log4j.Level
3
4 Logger.getLogger("org").setLevel(Level.ERROR)
5 Logger.getLogger("akka").setLevel(Level.ERROR)
6
```

Import the GraphX and RDD libraries.

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

**Step 2. Create a new dataset.**

**Define a simple list of vertices containing five international airports.**

Input:

```
1 val airports: RDD[(VertexId, String)] = sc.parallelize(
2   List((1L, "Los Angeles International Airport"),
3         (2L, "Narita International Airport"),
4         (3L, "Singapore Changi Airport"),
5         (4L, "Charles de Gaulle Airport"),
6         (5L, "Toronto Pearson International Airport"))
7
```

Output:

```
1 airports: org.apache.spark.rdd.RDD[(org.apache.spark.graphx.VertexId, String)] = ParallelCo
```

**Define a list of edges that will make up the flights.**

Two airports are connected in this graph if there is a flight between them. We will assign a made up flight number to each flight.

Input:

```
1 val flights: RDD[Edge[String]] = sc.parallelize(
2   List(Edge(1L,4L,"AA1123"),
3         Edge(2L, 4L, "3LS427"),
4         Edge(2L, 5L, "SQ9338"),
5         Edge(1L, 5L, "AA6653"),
6         Edge(3L, 4L, "SQ4521"))
7
```

Output:

```
1 flights: org.apache.spark.rdd.RDD[(org.apache.spark.graphx.Edge[String]] = ParallelCollectio
```

**Define the `flightGraph` graph from the airports vertices and the flights edges.**

Input:

```
1 val flightGraph = Graph(airports, flights)
2
```

Output:

```
1 flightGraph: org.apache.spark.graphx.Graph[String,String] = org.apache.spark.graphx.impl.Gr
```

**Print the departing and arrival airport and the flight number for each triplet in the `flightGraph` graph.**

Each triplet in the `flightGraph` graph represents a flight between two airports.

Input:

```
1 flightGraph.triplets.foreach(t => println("Departs from: " + t.srcAttr + " - Arrives at: " + t.dstAttr))
2
```

Output:

```
1 Departs from: Los Angeles International Airport - Arrives at: Charles de Gaulle Airport - Flight Num
2 Departs from: Narita International Airport - Arrives at: Charles de Gaulle Airport - Flight Num
3 Departs from: Singapore Changi Airport - Arrives at: Toronto Pearson International Airport - Flight Num
4 Departs from: Los Angeles International Airport - Arrives at: Toronto Pearson International
5 Departs from: Singapore Changi Airport - Arrives at: Charles de Gaulle Airport - Flight Num
```

**Define an `AirportInformation` class to store the airport city and code.**

Lets define a dataset with airport information so we can practice joining the airport information dataset with the datasets that we have already defined.

Input:

```
1 case class AirportInformation(city: String, code: String)
2
```

Output:

```
1 defined class AirportInformation
2
```

**Define the list of airport information vertices.**

Note: We do not have airport information defined for each airport in `flightGraph` and we have airport information for airports not in `flightGraph` graph.

Input:

```
1 val airportInformation: RDD[(VertexId, AirportInformation)] = sc.parallelize(
2   List((2L, AirportInformation("Tokyo", "NRT")),
3         (3L, AirportInformation("Singapore", "SIN")),
4         (4L, AirportInformation("Paris", "CDG")),
5         (5L, AirportInformation("Toronto", "YYZ")),
6         (6L, AirportInformation("London", "LHR")),
7         (7L, AirportInformation("Hong Kong", "HKG")))
8
```

Output:

```
1 airportInformation: org.apache.spark.rdd.RDD[(org.apache.spark.graphx.VertexId, AirportInfo
```

**Step 3. Join two datasets with `JoinVertices`.** In this first example we are going to use `joinVertices` to join the airport information `flightGraph` graph.

Create a mapping function that appends the city name to the name of the airport. The mapping function should return a string since that is the vertex attribute type of the `flightsGraph` graph.

Input:

```
1 def appendAirportInformation(id: VertexId, name: String, airportInformation: AirportInformation) =
2   name + " " + airportInformation.city
```

Output:

```
1 appendAirportInformation: (id: org.apache.spark.graphx.VertexId, name: String, airportInforma
```

Use `joinVertices` on `flightGraph` to join the `airportInformation` vertices to a new graph called `flightJoinedGraph` using the `appendAirportInformation` mapping function.

Input:

```
1 val flightJoinedGraph = flightGraph.joinVertices(airportInformation)(appendAirportInformation)
2 flightJoinedGraph.vertices.foreach(println)
3
```

Output:

```
1 (1,Los Angeles International Airport)
2 (4,Charles de Gaulle Airport:Paris)
3 (2,Narita International Airport:Tokyo)
4 (3,Singapore Changi Airport:Singapore)
5 (5,Toronto Pearson International Airport:Toronto)
6
```

**Step 4. Join two datasets with `outerJoinVertices`.** Use `outerJoinVertices` on `flightGraph` to join the `airportInformation` vertices with additional `airportInformation` such as city and code, to a new graph called `flightOuterJoinedGraph` using the `=>` operator which is just syntactic sugar for creating instances of functions.

Input:

```
1 val flightOuterJoinedGraph = flightGraph.outerJoinVertices(airportInformation)((_, name, airportInformation) => Airport(name, airportInformation.city, airportInformation.code))
2 flightOuterJoinedGraph.vertices.foreach(println)
3
```

Output:

```
1 ((1,(Los Angeles International Airport,None)))
2 ((4,(Charles de Gaulle Airport,Some(AirportInformation(Paris,CDG))))
3 ((2,(Narita International Airport,Some(AirportInformation(Tokyo,NRT))))
4 ((3,(Singapore Changi Airport,Some(AirportInformation(Singapore,SIN))))
5 ((5,(Toronto Pearson International Airport,Some(AirportInformation(Toronto,YYZ))))
```

Use `outerJoinVertices` on `flightGraph` to join the `airportInformation` vertices with additional `airportInformation` such as city and code, to a new graph called `flightOuterJoinedGraphTwo` but this time printing `'NA'` if there is no additional information.

Input:

```
1 val flightOuterJoinedGraphTwo = flightGraph.outerJoinVertices(airportInformation)((_, name, airportInformation) => Airport(name, airportInformation.city, airportInformation.code))
2 flightOuterJoinedGraphTwo.vertices.foreach(println)
3
```

Output:

```
1 ((1,(Los Angeles International Airport,NA)))
2 ((4,(Charles de Gaulle Airport,Some(AirportInformation(Paris,CDG))))
3 ((2,(Narita International Airport,Some(AirportInformation(Tokyo,NRT))))
4 ((3,(Singapore Changi Airport,Some(AirportInformation(Singapore,SIN))))
5 ((5,(Toronto Pearson International Airport,Some(AirportInformation(Toronto,YYZ))))
```

**Step 5. Create a new return type for the joined vertices.** Create a case class called `Airport` to store the information for the name, city, and code of the airport.

Input:

```
1 case class Airport(name: String, city: String, code: String)
2
```

Output:

```
2
```

Print the `airportInformation` with the name, city, and code within each other.

Input:

```
1 val flightOuterJoinedGraphThree = flightGraph.outerJoinVertices(airportInformation)((_, name, airportInformation) => Airport(name, airportInformation.city, airportInformation.code))
2 flightOuterJoinedGraphThree.vertices.foreach(println)
3
```

Output:

```
1 ((1,Airport(Los Angeles International Airport,NA)))
2 ((3,Airport(Singapore Changi Airport,Singapore,SIN))
3 ((2,Airport(Narita International Airport,Tokyo,NRT)))
4 ((4,Airport(Charles de Gaulle Airport,Paris,CDG)))
5 ((5,Airport(Toronto Pearson International Airport,Toronto,YYZ)))
```

**Go to next item** **Completed**