

# Dig Data Management Assignment 8

信科四 郑元嘉 1800920541

In this report, we'll implement 'PageRank' by GraphX based on Spark.

Get it straight, we simply assign score of 1 to each nodes initially, so the sum of all nodes' scores is not equal to 1.

We will run `scala` script in `spark-shell` .

As of the network data, it's converted into structured `.txt` file from and `.png` file:

```
1  A C F D B
2  B A D G E
3  C A D F
4  D A B C E F G
5  E B D G
6  F A C D G H
7  G B D E F H
8  H F G I
9  I H J
10 J I
```

Take the 1st row as an instance, 'A C F D B' indicates node `A` is connected to nodes `C` , `F` , `D` and `B` .

## Scala Script

The script is going to be broken down into segments.

```
1  // val conf = new SparkConf().setAppName("pageRank")
2  // val sc = new SparkContext(conf)
3
4  val fp = "/media/jack/File/hadoopHw/hw8/jack/network.txt"
5
6  var mapping = sc.textFile(fp)
7  .flatMap(line => line.split(" "))
8  .distinct
9  .map(node => (MurmurHash3.stringHash(node).toLong, node))
```

First things first. I'm confronted by my first bug when I'm trying to initiate `SparkContext` because the spark shell has already initiated one and the default setting allows only one `SparkContext` instance at a time.

Function `sc.textFile()` is applied here to read `.txt` file to build up one table mapping the hashed id number back into the original character name (such as A, B, ...). It returns `RDD` object instead of any related objects of `collection`, so we pay more attention when we coping with type transformation.

We use `flatMap` to extract all possible characters, use `distinct` to capture distinct ones, and eventually use `map` to create `Tuple` pairs in the favor of `MurmurHash3`'s hashing.

```
1 var links = List[Tuple2[Long, Long]]()
2
3 for (line <- Source.fromFile(fp).getLines) {
4     var s = line.split(" ")(0)
5     line.split(" ").drop(1).map(node => links = (MurmurHash3.stringHash(s).toLo
6 }
7
8 var linksRdd = sc.parallelize(links).persist()
```

Alternatively, we apply function `Source.fromFile()` to read and parse the `.txt` file into sourceNode-destinationNode pairs which follow the rules for the input of the function `Graph.fromEdgeTuples()` shown later. Note that due to object `Graph.Vertice`, the type of the vertices (nodes) is required as `Long` rather than default `Int` or else.

A common Trick of list is played here, we need to use `::` to add new element to the list (what's more, only add at the 'head' can be achieved), and be careful of type assignment for new empty `List` object.

To transform the pairs `List` to a `RDD`, we apply `parallelize()`. What's more, call `persist()` function to store it in RAM to speed up the task.

```
1 val graph = Graph.fromEdgeTuples( linksRdd, 1 )
2
3 val ranks =
4     graph
5     .pageRank( 0.01 )
6     .vertices
7
8 val fullMap =
9     ranks
10    .join( mapping )
11    .map( row => row._2)
12    .sortBy(_._2)
13    .collect
14
```

```
.foreach(println)
}
```

Function `Graph.fromEdgeTuples` is one of the ways to create one `Graph` object. We assign `1` as all vertices' attributes.

Then, we call function `pageRank` to run our page ranking task.

Lastly, we use `join()` and table `mapping` to map the hashed ids back the original node character names, and use `map()` and `sortBy()` to get a clear and wanted form of data.

Here's a point that should be kept in mind. Call `foreach(println)` to show the result after calling `collect` or other action operations; otherwise you'll get unexpected results (like unsorted results).

## Result

```
scala> :load /media/jack/File/hadoopHw/hw8/jack/pageRank.scala
Loading /media/jack/File/hadoopHw/hw8/jack/pageRank.scala...
import scala.io.Source
import org.apache.spark.graphx.Graph
import scala.util.hashing.MurmurHash3
defined object PageRank
(1.024373490355488,A)
(1.024373490355488,B)
(0.8016717537932411,C)
(1.472866894710858,D)
(0.8016717537932411,E)
(1.2832802325323818,F)
(1.2832802325323818,G)
(0.9488863714440244,H)
(0.8446560052929347,I)
(0.5149397751899604,J)
```