# CES Data Scientist – PageRank Exam

Youcef KACER (`youcef.kacer@telecom-paristech.fr`)

28 septembre 2016

**Exercise 1 :** We expose here after the system of linear equations that solves PageRank scores associated to graph $G$ ($r_i$ being PageRank score for page $i$) :

$$
\begin{aligned}
r_1 &= r_4 \\
r_2 &= 1/2 * r_1 \\
r_3 &= 1/2 * r_1 + r_2 \\
r_4 &= r_3 \\
1 &= r_1 + r_2 + r_3 + r_4
\end{aligned}
$$

We can solve manually the system :

$$
\begin{aligned}
r_1 &= r_3 \\
r_1 &= r_4 \\
r_2 &= 1/2 * r_1 \\
1 &= 7/2 * r_1
\end{aligned}
$$

$$
\begin{aligned}
r_1 &= 2/7 \\
r_2 &= 1/7 \\
r_3 &= 2/7 \\
r_4 &= 2/7
\end{aligned}
$$

Otherwise, we can solve the system using previous PageRank lab using MapReduce. We just need to adapt *edge_list.txt* to the graph $G$. In this file, each line represents a page indice and the page indices it links to :

```
1 2 3
2 3
3 4
4 1
```

We have run our python map/reduce implementation of PageRank. We use a teleport coefficient of 0.75, and an error criteria of 0.01. Figure 1 shows correspoding standard output with PageRank scores results at the end.
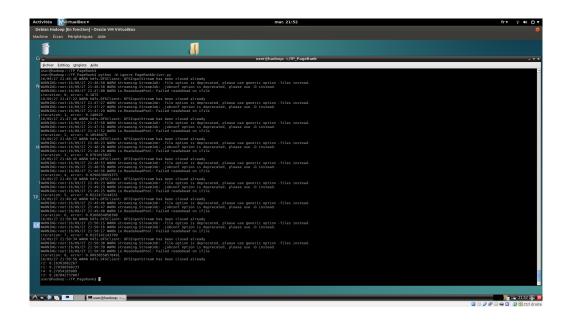
FIGURE 1 – Python map/reduce for PageRank scores associated to graph $G$

**Exercise 2 :** Suppose we have a file containing integers with theirs indices as follow :

```
i n
0 256
1 35
2 4122
3 96
...
```

We present here after a pseudo-code for map/reduce function to find maximum value of this list of integers :

```
map(key=i,value=n)
{
  return 0,n
}

reduce(key,values) // values is an array containing all the integers because
{                  // map has emitted each of them with the same key (0).
  MAX = MIN_POSSIBLE_VALUE
  for n in values
  {
    if n>max
    {
      max = n
    }
  }
  return 0,max
}
```

**Exercise 3 :**

**Exercise 4 :**