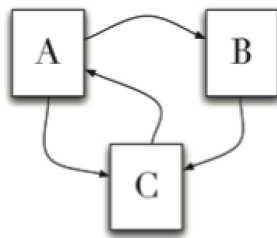# page_rank.py

In this code I have implemented the google's page rank algorithm stated below. Please ignore any mistakes in the usage of python or the code as it was my first attempt at python.

# PageRank



- PageRank (*PR*) of page C = *PR*(A)/2 + *PR*(B)/1
- More generally,

$$PR(u) = \sum_{v \in B_u} \frac{PR(v)}{L_v}$$

  — where $B_u$ is the set of pages that point to *u*, and $L_v$ is the number of outgoing links from page *v* (not counting duplicate links)

## Algorithm:

```
 1: procedure PAGERANK(G)
 2:        ▷ G is the web graph, consisting of vertices (pages) and edges (links).
 3:     (P, L) ← G                              ▷ Split graph into pages and links
 4:     I ← a vector of length |P|              ▷ The current PageRank estimate
 5:     R ← a vector of length |P|     ▷ The resulting better PageRank estimate
 6:     for all entries Iᵢ ∈ I do
 7:         Iᵢ ← 1/|P|                     ▷ Start with each page being equally likely
 8:     end for
 9:     while R has not converged do
10:         for all entries Rᵢ ∈ R do
11:             Rᵢ ← λ/|P| ▷ Each page has a λ/|P| chance of random selection
12:         end for
13:         for all pages p ∈ P do
14:             Q ← the set of pages such that (p, q) ∈ L and q ∈ P
15:             if |Q| > 0 then
16:                 for all pages q ∈ Q do
17:                     R_q ← R_q + (1 − λ)I_p/|Q|        ▷ Probability I_p of being at
        page p
18:                 end for
19:             else
20:                 for all pages q ∈ P do
21:                     R_q ← R_q + (1 − λ)I_p/|P|
22:                 end for
23:             end if
24:             I ← R                         ▷ Update our current PageRank estimate
25:         end for
26:     end while
27:     return R
28: end procedure
```

# A PageRank Implementation

Iteration:

- Steps:
  1. Make a new output file, R.
  2. Read L and I in parallel (since they're all sorted by URL).
  3. For each unique source URL, determine whether it has any outgoing links:
  4. If not, add its current PageRank value to the sum: T (terminals).
  5. If it does have outgoing links, write (source_url, dest_url, Ip/|Q|), where Ip is the current PageRank value, |Q| is the number of outgoing links, and dest_url is a link destination.
     Do this for all outgoing links.  Write this to R.
  6. Sort R by destination URL.
  7. Scan R and I at the same time.  The new value of Rp is:
     (1 - lambda) / #D (a fraction of the sum of all pages)
     plus: lambda * sum(T) / #D (the total effect from terminal pages),
     plus: lambda * all incoming mass from step 5. ()
  8. Check for convergence
  9. Write new Rp values to a new I file.