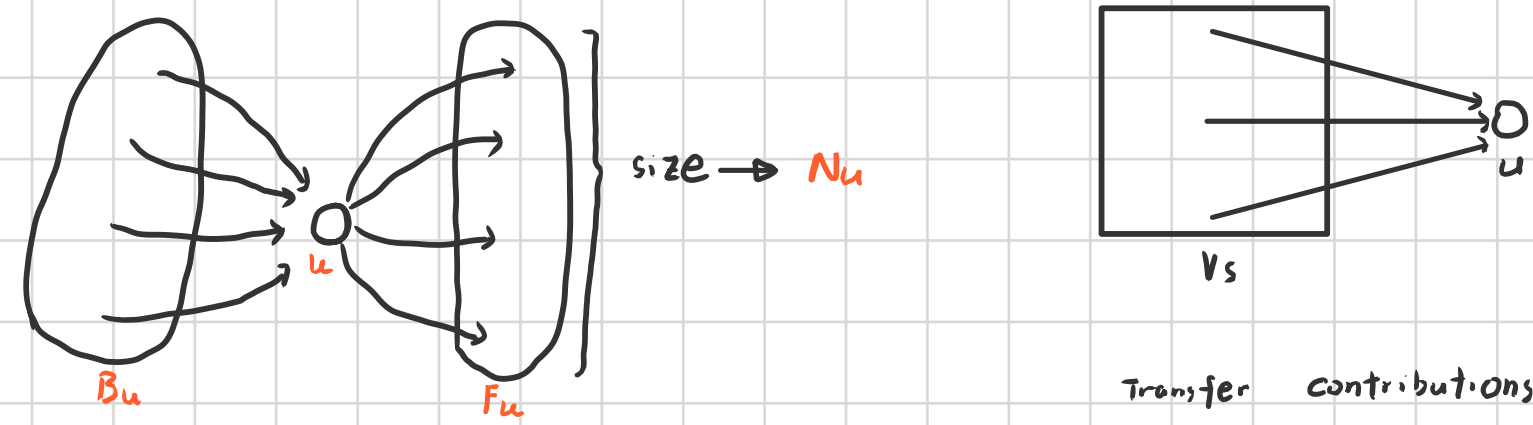


• PageRank Algorithm's basic theorems

[Boy highlight 01]

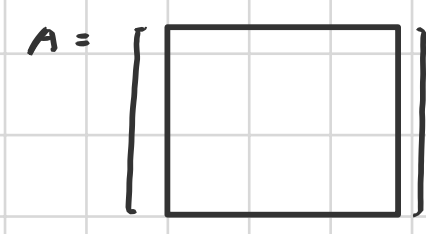


$$R(u) = C \sum_{v \in B_u} \frac{R(v)}{N_v}$$

C < 1: there are a number of pages with on forward links, and their weight is lost from the system

[Boy highlight 02] Use matrix for computing & Simple Model

Initial Status:



$$A_{u,v} = 1/N_u \text{ if } u \rightarrow v$$

$$A_{u,v} = 0 \text{ if } u \nrightarrow v$$

Aim: find a vector $R = (r_1, r_2, \dots, r_n)^T$, $R[u]$ = Page i's important value
then

$$C A^T R = C \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} = \begin{bmatrix} C \sum_{j \in B_u} \frac{R_j}{N_j} \\ \vdots \end{bmatrix}$$

$$A^T R = \frac{1}{C} R$$

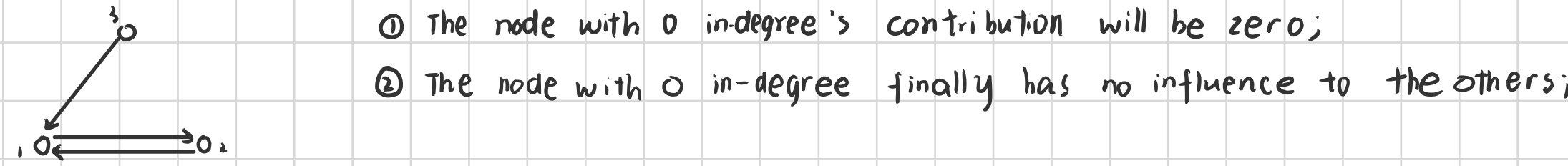
eigenvector
want dominant eigenvector

$$A^T_{u,v} = 1/N_u \text{ if } v \rightarrow u$$

$$A^T_{u,v} = 0 \text{ if } v \nrightarrow u$$

[Boy highlight 02 below] Two problems & Completed / Random Surfer Model

"rank sink": there are some nodes with 0 in-degree



$$A = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix}$$

for each $A_{u,v}$: $u \rightarrow v$
row vector = out-edges

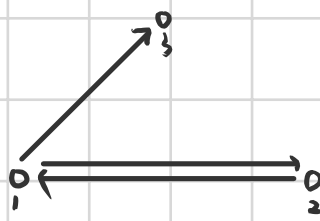
$$A^T = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

for each $A^T_{u,v}$: $v \rightarrow u$
row vector = in-edges

① so if there is a node with 0 in-degree, surely, the eigenvector's corresponding value will be zero
② because the 0 in-degree's contribution is zero, so its out-edge will have no influence to other node

"rank leak": there are some nodes with 0-outdegree

① Finally, all node which link to the node with 0-outdegree will be zero contribution



$$A = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

for each A_{ij} : $i \rightarrow j$
row vector = out-edges

$$A^T = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix}$$

for each A^T_{ij} : $j \rightarrow i$
row vector = in-edges
always < 1

For each time, the node with 0-outdegree have on contribution to the others, but the others which link to it save < 1 ratio than last time, so finally, their contribution $\rightarrow 0$

Solution:

Let the "pure matrix A" becomes a "stable probability stochastic matrix A"

$$R'(u) = C \sum_{v \in B_u} \frac{R'(v)}{N_v} + C E(u)$$

$$R' = C(A^T + E \times 1)R'$$

$$R(u) = \frac{1-C}{N} + C \sum_{v \in B_u} \frac{R(v)}{N_v}$$

All pages Number
Damping Factor (usually is 0.85)
Random Surfer Model

Proof of convergence: Markov chain convergence theorem

• Implement PageRank Algorithm by Hadoop MapReduce

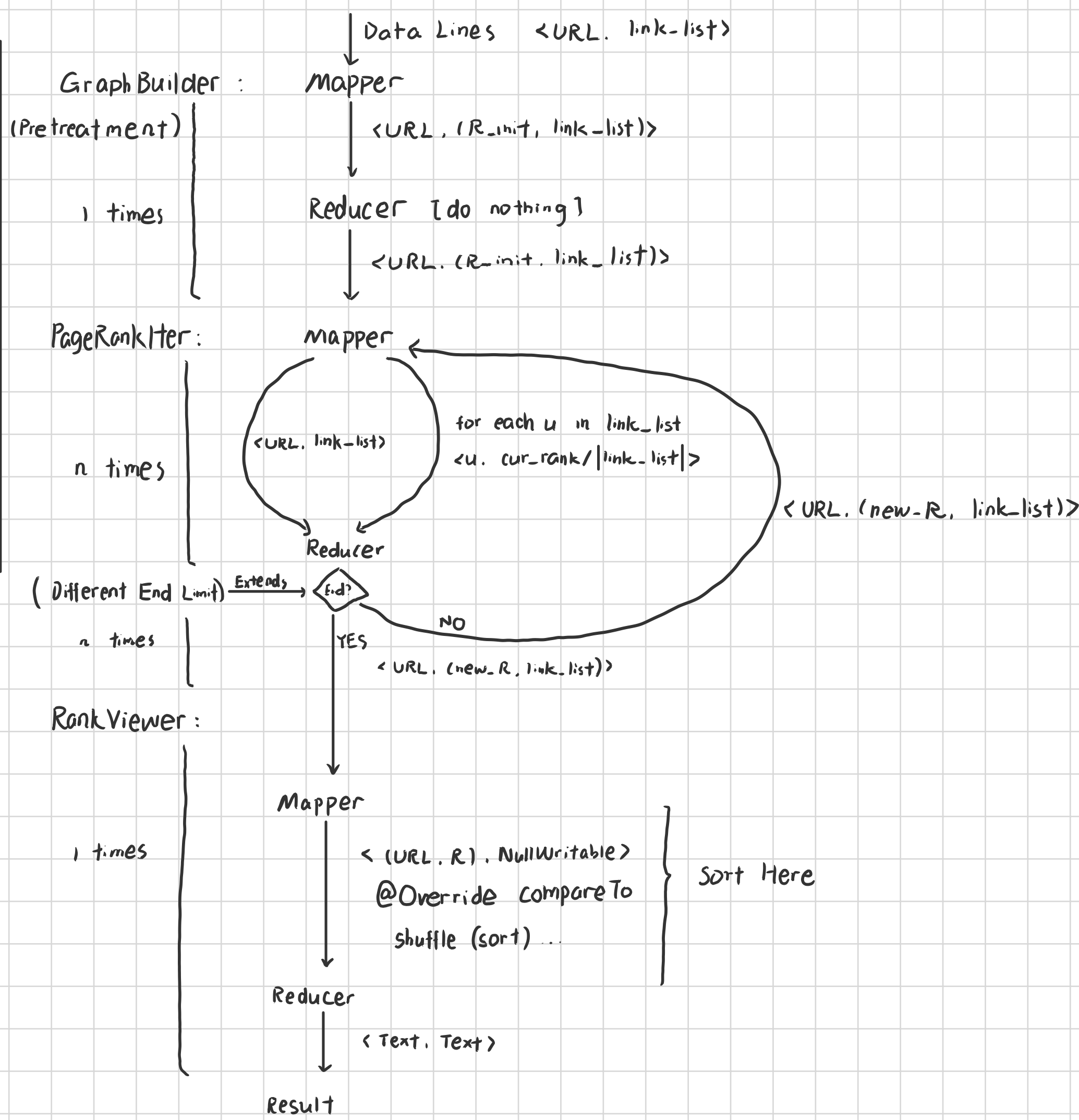
Implement Idea:

- 3 parts:
1. GraphBuilder
 2. PageRankIter
 3. RankViewer

2 parallel parts:

$$R(u) = \frac{1-d}{N} + d \sum_{v \in B_u} \frac{R(v)}{N_v}$$

especially / 1. Mapper parallel op
2. Reducer parallel op



Hard Points

- ① How to deal with input/output files store, especially for "Pretreatment" \rightarrow "Real PageRank" ☒
- ② How to use Hadoop MapReduce iteration workflow? ☒
- ③ How to implement and apply a extendable end limit class & methods? ☐

More problem: How to set initial PageRank value for each page? ☐
How to get all pages number which is the parameter N in $R(u) = \frac{1-d}{N} + d \sum_{v \in B_u} \frac{R(v)}{N_v}$ ☒
(use a single process for scanning, but how to optimize it?) ☐