

## 1.4.1 Computing PageRank

### □ Key step is matrix-vector multiplication

➤  $\mathbf{r}^{\text{new}} = \mathbf{A} \cdot \mathbf{r}^{\text{old}}$

### □ Easy if we have enough main memory to hold $\mathbf{A}$ , $\mathbf{r}^{\text{old}}$ , $\mathbf{r}^{\text{new}}$

### □ Say $N = 1$ billion (十亿) pages

- We need 4 bytes for each entry (say)
- $\mathbf{r}^{\text{old}}$ ,  $\mathbf{r}^{\text{new}}$ : 2 billion entries for vectors, approx 8GB
- **A: Matrix  $\mathbf{A}$  has  $N^2$  entries**
  - $10^{18}$  is a large number!

$$\mathbf{A} = \beta \cdot \mathbf{M} + (1 - \beta) [1/N]_{N \times N}$$

$$\mathbf{A} = 0.8 \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & 0 & 0 \\ 0 & \frac{1}{2} & 1 \end{bmatrix} + 0.2 \begin{bmatrix} 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \end{bmatrix}$$

$$= \begin{bmatrix} 7/15 & 7/15 & 1/15 \\ 7/15 & 1/15 & 1/15 \\ 1/15 & 7/15 & 13/15 \end{bmatrix}$$

【备注】观察矩阵 $\mathbf{A}$ 和 $\mathbf{M}$ 可知,  
Matrix  $\mathbf{A}$  **dense matrix!**  
Matrix  $\mathbf{M}$  **sparse matrix**

## 1.4.1 Rearranging the Equation

$$\square \mathbf{r} = \mathbf{A} \cdot \mathbf{r}, \text{ where } A_{ji} = \beta M_{ji} + \frac{1-\beta}{N}$$

$$\square r_j = \sum_{i=1}^N A_{ji} \cdot r_i$$

$$\begin{aligned} \square r_j &= \sum_{i=1}^N \left[ \beta M_{ji} + \frac{1-\beta}{N} \right] \cdot r_i \\ &= \sum_{i=1}^N \beta M_{ji} \cdot r_i + \frac{1-\beta}{N} \sum_{i=1}^N r_i \\ &= \sum_{i=1}^N \beta M_{ji} \cdot r_i + \frac{1-\beta}{N} \end{aligned}$$

since  $\sum r_i = 1$

$$\square \text{So we get: } \mathbf{r} = \beta \mathbf{M} \cdot \mathbf{r} + \left[ \frac{1-\beta}{N} \right]_N$$

$[x]_N \dots$  a vector of length  $N$  with all entries  $x$

$$\mathbf{A} = \beta \mathbf{M} + (1-\beta) \left[ \frac{1}{N} \right]_{N \times N}$$

**Note:** Here we assumed  $\mathbf{M}$  has no dead-ends

## 1.4.1 Sparse Matrix Formulation


□ We just rearranged the **PageRank equation**

$$r = \beta M \cdot r + \left[ \frac{1 - \beta}{N} \right]_N$$

where  $[(1-\beta)/N]_N$  is a vector with all  $N$  entries  $(1-\beta)/N$

□  $M$  is a **sparse matrix!** (with **no dead-ends**)

➤  $N$  nodes, 10 links per node, approx  $10N$  entries

$o(2N+N*N)$    
 $o(2N+10N)$

□ So in each iteration, we need to:

➤ Compute  $r^{new} = \beta M \cdot r^{old}$

➤ Add a constant value  $(1-\beta)/N$  to each entry in  $r^{new}$

- Note if  $M$  contains dead-ends then  $\sum_j r_j^{new} < 1$ , and we have to renormalize  $r^{new}$  so that it sums to 1

# 1.4.1 The Complete PageRank Algorithm

## Input: Graph $G$ and parameter $\beta$

- Directed graph  $G$  (can have **spider traps** and **dead ends**)
- Parameter  $\beta$

$$r = \beta M \cdot r + \left[ \frac{1 - \beta}{N} \right]_N$$

## Output: PageRank vector $r^{new}$

- **Set:**  $r_j^{old} = \frac{1}{N}$
- **repeat until convergence:**  $\sum_j |r_j^{new} - r_j^{old}| < \varepsilon$ 
  - $\forall j: r_j'^{new} = \sum_{i \rightarrow j} \beta \frac{r_i^{old}}{d_i}$   
 $r_j'^{new} = 0$  if in-degree of  $j$  is 0
  - **Now re-insert the leaked PageRank:**  
 $\forall j: r_j^{new} = r_j'^{new} + \frac{1 - S}{N}$  where:  $S = \sum_j r_j'^{new}$
  - $r^{old} = r^{new}$

- If the graph has no dead-ends then the amount of leaked PageRank is **1- $\beta$** .
- But since we have **dead-ends**, the amount of leaked PageRank may **be larger**.
- Hence, we have to explicitly account for it by computing **S**.

## 1.4.2 Sparse Matrix Encoding

□ Encode sparse matrix  $M$  using only nonzero entries

source node	degree	destination nodes
0	3	1, 5, 7
1	5	17, 64, 113, 117, 245
2	2	13, 23

- Space proportional roughly to **number of links**
- Say  $10N$  ( $N$  nodes, 10 links per node), or  $4 \times 10^9$  billion = **40GB**, e.g.  $N = 1$  billion (十亿)
- $M$  still won't fit in memory, but will **fit on disk**

## 1.4.2 Basic Update Algorithm

Assume enough RAM to fit  $r^{new}$  into memory. Store  $r^{old}$  and matrix  $M$  on disk

1 step of power-iteration is:

Initialize all entries of  $r^{new} = (1-\beta) / N$

For each page  $i$  (of out-degree  $d_i$ ):

Read into memory:  $i, d_i, dest_1, \dots, dest_{d_i}, r^{old}(i)$

For  $j = 1 \dots d_i$

$r^{new}(dest_j) += \beta r^{old}(i) / d_i$

0	
1	
2	
3	
4	
5	
6	

$r^{new}$

source degree destination

0	3	1, 5, 6
1	4	17, 64, 113, 117
2	2	13, 23

For one iteration

$r^{old}$

0	
1	
2	
3	
4	
5	
6	

For one iteration

## 1.4.2 Analysis of Basic Update

### □ Assume enough RAM to fit $r^{new}$ into memory

- Store  $r^{old}$  and matrix  $M$  on disk

### □ In each iteration, we have to:

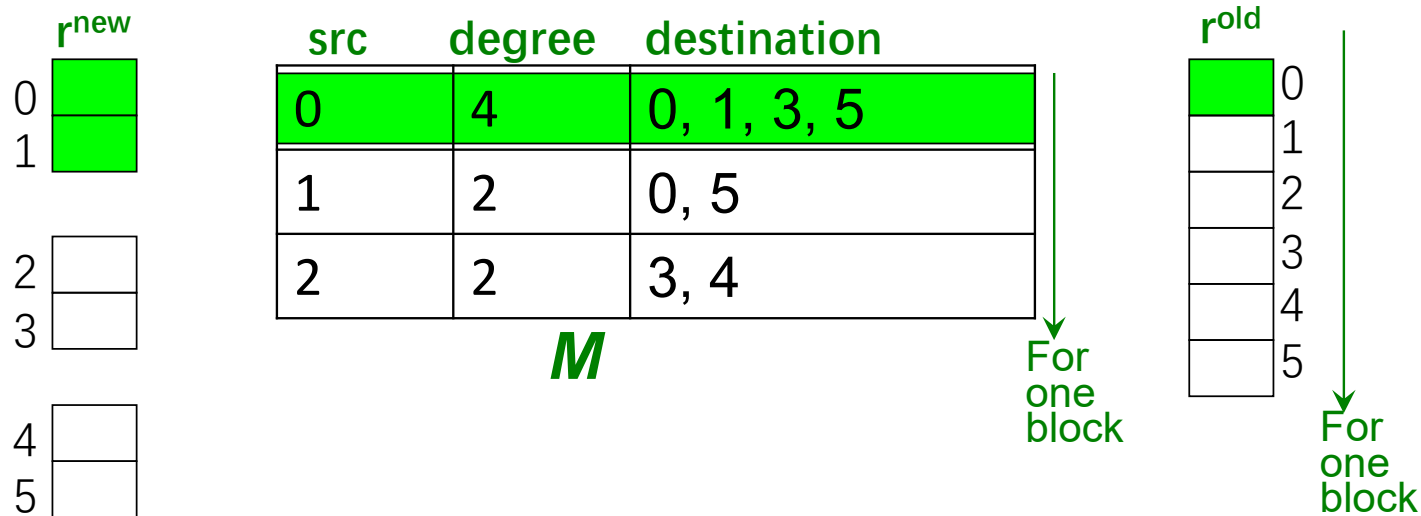
- Read  $r^{old}$  and  $M$
- Write  $r^{new}$  back to disk
- **Cost per iteration of Power method:**  
 $= 2|r| + |M|$

### □ Question:

- What if we could not even fit  $r^{new}$  in memory?

## 1.4.3 Block-based Update Algorithm

- Break  $r^{\text{new}}$  into  $k$  blocks that fit in memory



- Scan  $M$  and  $r^{\text{old}}$  once for each block
- In each iteration total  $k$  blocks, then  $k$  scans  $M$  and  $r^{\text{old}}$



## 1.4.3 Analysis of Block Update

### ❑ Similar to nested-loop join in databases

- Break  $r^{\text{new}}$  into  $k$  blocks that fit in memory
- Scan  $M$  and  $r^{\text{old}}$  once for each block

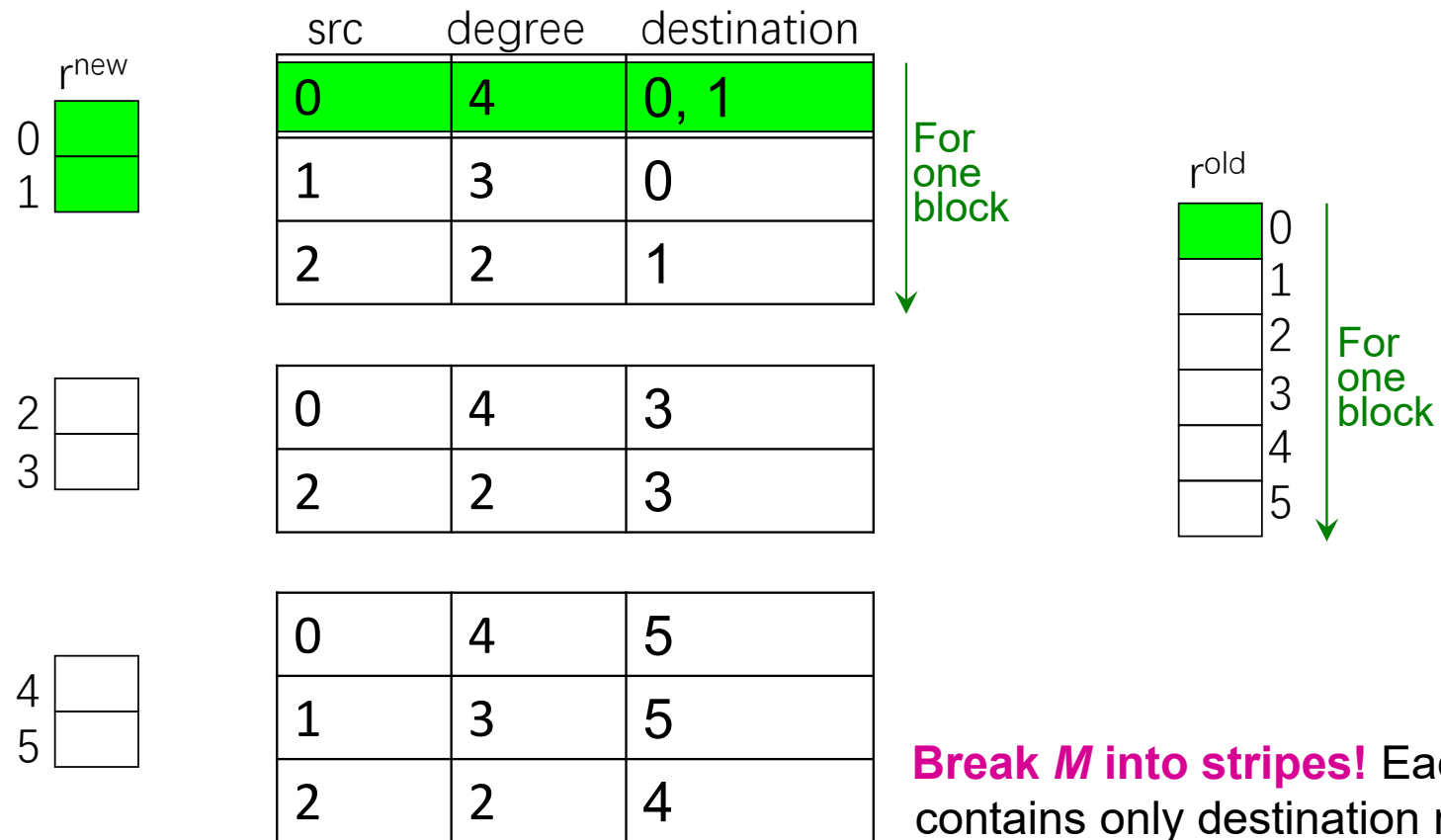
### ❑ Total cost:

- $k$  scans of  $M$  and  $r^{\text{old}}$
- Write  $r^{\text{new}}$  back to disk ( $k$  blocks)
- **Cost per iteration of Power method:**  
 $k(|M| + |r|) + |r| = k|M| + (k+1)|r|$

### ❑ Can we do better?

- **Hint:**  $M$  is much bigger than  $r$  (approx 10-20x), so we must **avoid reading it  $k$  times** per iteration

## 1.4.4 Block-Stripe Update Algorithm



**Break  $M$  into stripes!** Each stripe contains only destination nodes in the corresponding block of  $r^{\text{new}}$

## 1.4.4 Analysis of Block-Stripe Update

### □ Break $M$ into stripes

➤ Each stripe contains only destination nodes in the corresponding block of  $r^{\text{new}}$

### □ Some additional overhead per stripe

➤ But it is usually worth it

### □ Cost per iteration of Power method:

$$= |M|(1+\varepsilon) + (k+1)|r|$$

# Some Problems with PageRank

## ❑ Measures generic popularity of a page

- Biased against topic-specific authorities
- **Solution:** Topic-Specific PageRank (**next**)

## ❑ Susceptible to Link spam

- Artificial link topographies created in order to boost page rank
- **Solution:** TrustRank (**next**)

## ❑ Uses a single measure of importance

- Other models of importance
- **Solution:** Hubs-and-Authorities (**next**)