

How do we really compute the PageRank?

Mining of Massive Datasets
Leskovec, Rajaraman, and Ullman
Stanford University



Computing Page Rank

- **Key step is matrix-vector multiplication**

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

- Easy if we have enough main memory to hold A , r^{old} , r^{new}

- **Say $N = 1$ billion pages**

- We need 4 bytes for each entry (say)

- 2 billion entries for vectors, approx 8GB

- **Matrix A has N^2 entries**

- 10^{18} is a large number!

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

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

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

Matrix Formulation

- Suppose there are N pages
- Consider page j , with d_j out-links
- We have $M_{ij} = 1/d_j$ when $j \rightarrow i$
and $M_{ij} = 0$ otherwise
- **The random teleport is equivalent to:**
 - Adding a **teleport link** from j to every other page and setting transition probability to $(1-\beta)/N$
 - Reducing the probability of following each out-link from $1/d_j$ to β/d_j
 - **Equivalent:** Tax each page a fraction $(1-\beta)$ of its score and redistribute evenly

Rearranging the Equation

- $\mathbf{r} = \mathbf{A} \cdot \mathbf{r}$, where $A_{ij} = \beta M_{ij} + \frac{1-\beta}{N}$
- $r_i = \sum_{j=1}^N A_{ij} \cdot r_j$
- $r_i = \sum_{j=1}^N \left[\beta M_{ij} + \frac{1-\beta}{N} \right] \cdot r_j$
 $= \sum_{j=1}^N \beta M_{ij} \cdot r_j + \sum_{j=1}^N \frac{1-\beta}{N} r_j$
 $= \sum_{j=1}^N \beta M_{ij} \cdot r_j + \frac{1-\beta}{N}$ since $\sum r_j = 1$
- So we get: $\mathbf{r} = \beta \mathbf{M} \cdot \mathbf{r} + \left[\frac{1-\beta}{N} \right]_N$

此处把M提取出来，因为M是稀疏的

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

$[\mathbf{x}]_N$... a vector of length N with all entries \mathbf{x}

Sparse Matrix Formulation

- We just rearranged the **PageRank equation**

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

- where $[(1-\beta)/N]_N$ is a vector with all N entries $(1-\beta)/N$
- \mathbf{M} is a **sparse matrix!** (with no dead-ends)
 - 10 links per node, approx $10N$ entries
- **So in each iteration, we need to:**
 - Compute $\mathbf{r}^{\text{new}} = \beta \mathbf{M} \cdot \mathbf{r}^{\text{old}}$
 - Add a constant value $(1-\beta)/N$ to each entry in \mathbf{r}^{new}
 - **Note if \mathbf{M} contains dead-ends then $\sum_i r_i^{\text{new}} < 1$ and we also have to renormalize \mathbf{r}^{new} so that it sums to 1**

PageRank: The Complete Algorithm

- Input: Graph G and parameter β

- Directed graph G with spider traps and dead ends
- Parameter β

- Output: PageRank vector r

- **Set:** $r_j^{(0)} = \frac{1}{N}, \quad t = 1$

- **do:**

- $\forall j: r_j'^{(t)} = \sum_{i \rightarrow j} \beta \frac{r_i^{(t-1)}}{d_i}$

- $r_j'^{(t)} = 0$ if in-deg. of j is 0

Dead Ends

- **Now re-insert the leaked PageRank:**

- $\forall j: r_j^{(t)} = r_j'^{(t)} + \frac{1-S}{N}$ where: $S = \sum_j r_j'^{(t)}$

- $t = t + 1$

- **while** $\sum_j |r_j^{(t)} - r_j^{(t-1)}| > \varepsilon$