# COMP 345: Data Mining
# More on PageRank

Slides Adapted From: www.mmds.org (Mining Massive Datasets)

Rhodes College

---

# Announcements

- Assignment 6
  - due Wed. Nov. 14th/Thurs. Nov. 15th

2

# MapReduce Quiz Problem

Suppose our input data to a map-reduce operation consists of integer values (the keys are not important). The map function takes an integer i and produces the list of pairs (p,i) such that p is a prime divisor of i. For example, map(12) = [(2,12), (3,12)]. The reduce function is addition. That is, reduce(p, [i , i , ...,i ]) is (p,i +i +...+i ). Compute the output, if the input is the set of integers 15, 21, 24, 30, 49. Then, identify, in the list below, one of the pairs in the output.

a.   (7, 70)
b.   (5, 49)
c.   (2, 47)
d.   (6, 54)

3

# PageRank:
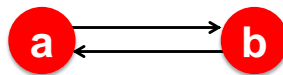# The Google Formulation

# PageRank: Three Questions

$$r_j^{(t+1)} = \sum_{i \to j} \frac{r_i^{(t)}}{d_i} \quad \text{or equivalently} \quad r = Mr$$

- **Does this converge?**

- **Does it converge to what we want?**

- **Are results reasonable?**

5

# Does this converge?

a ⇄ b
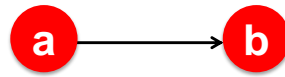
$$r_j^{(t+1)} = \sum_{i \to j} \frac{r_i^{(t)}}{d_i}$$

- **Example:**

$$\begin{matrix} r_a \\ r_b \end{matrix} = \begin{matrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{matrix}$$

Iteration 0, 1, 2, …

6

# Does it converge to what we want?

$$a \longrightarrow b \qquad r_j^{(t+1)} = \sum_{i \to j} \frac{r_i^{(t)}}{d_i}$$

- **Example:**

$$\begin{matrix} r_a \\ r_b \end{matrix} = \begin{matrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{matrix}$$

Iteration 0, 1, 2, …

7

# PageRank: Problems

Dead end

**2 problems:**
- **(1)** Some pages are
  **dead ends** (have no out-links)
  - Random walk has "nowhere" to go to
  - Such pages cause importance to "leak out"

  Spider trap

- **(2) Spider traps:**
  (all out-links are within the group)
  - Random walked gets "stuck" in a trap
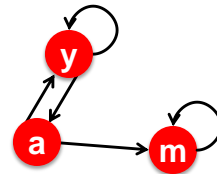  - And eventually spider traps absorb all importance

8

# Problem: Spider Traps

- **Power Iteration:**
  - Set $r_j = 1$
  - $r_j = \sum_{i \to j} \frac{r_i}{d_i}$
    - And iterate



|   | y | a | m |
|---|---|---|---|
| y | ½ | ½ | 0 |
| a | ½ | 0 | 0 |
| m | 0 | ½ | 1 |

m is a spider trap

$r_y = r_y/2 + r_a/2$
$r_a = r_y/2$
$r_m = r_a/2 + r_m$

- **Example:**

$$\begin{bmatrix} r_y \\ r_a \\ r_m \end{bmatrix} = \begin{array}{ccccc} 1/3 & 2/6 & 3/12 & 5/24 & \quad 0 \\ 1/3 & 1/6 & 2/12 & 3/24 \quad \ldots & \quad 0 \\ 1/3 & 3/6 & 7/12 & 16/24 & \quad 1 \end{array}$$

Iteration 0, 1, 2, …

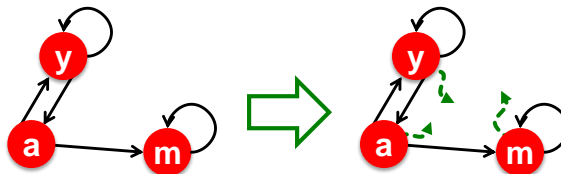All the PageRank score gets "trapped" in node m.

9

# Solution: Teleports!

- **The Google solution for spider traps: At each time step, the random surfer has two options**
  - With prob. $\beta$, follow a link at random
  - With prob. **1-$\beta$**, jump to some random page
  - Common values for $\beta$ are in the range 0.8 to 0.9
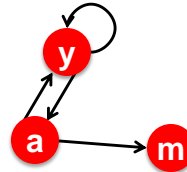- **Surfer will teleport out of spider trap within a few time steps**

10

5

# Problem: Dead Ends

- **Power Iteration:**
  - Set $r_j = 1$
  - $r_j = \sum_{i \to j} \frac{r_i}{d_i}$
    - And iterate

|   | y | a | m |
|---|---|---|---|
| y | ½ | ½ | 0 |
| a | ½ | 0 | 0 |
| m | 0 | ½ | 0 |

$r_y = r_y/2 + r_a/2$
$r_a = r_y/2$
$r_m = r_a/2$

- **Example:**

$$\begin{bmatrix} r_y \\ r_a \\ r_m \end{bmatrix} = \begin{matrix} 1/3 & 2/6 & 3/12 & 5/24 & & 0 \\ 1/3 & 1/6 & 2/12 & 3/24 & \dots & 0 \\ 1/3 & 1/6 & 1/12 & 2/24 & & 0 \end{matrix}$$
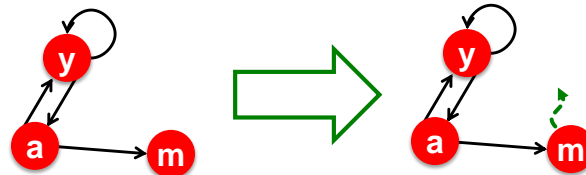
Iteration 0, 1, 2, …

Here the PageRank "leaks" out since the matrix is not stochastic.

# Solution: Always Teleport!

- **Teleports:** Follow random teleport links with probability 1.0 from dead-ends
  - Adjust matrix accordingly



|   | y | a | m |
|---|---|---|---|
| y | ½ | ½ | 0 |
| a | ½ | 0 | 0 |
| m | 0 | ½ | 0 |

|   | y | a | m |
|---|---|---|---|
| y | ½ | ½ | ⅓ |
| a | ½ | 0 | ⅓ |
| m | 0 | ½ | ⅓ |

# Why Teleports Solve the Problem?

**Why are dead-ends and spider traps a problem and why do teleports solve the problem?**

- **Spider-traps** are not a problem, but with traps PageRank scores are **not** what we want
  - **Solution:** Never get stuck in a spider trap by teleporting out of it in a finite number of steps
- **Dead-ends** are a problem
  - The matrix is not column stochastic so our initial assumptions are not met
  - **Solution:** Make matrix column stochastic by always teleporting when there is nowhere else to go

# Solution: Random Teleports

- **Google's solution that does it all:**
  At each step, random surfer has two options:
  - With probability $\beta$, follow a link at random
  - With probability $1\text{-}\beta$, jump to some random page

- **PageRank equation** [Brin-Page, 98]

$$r_j = \sum_{i \to j} \beta \, \frac{r_i}{d_i} + (1 - \beta)\frac{1}{N}$$

$d_i$ ... out-degree of node i

This formulation assumes that $M$ has no dead ends. We can either preprocess matrix $M$ to remove all dead ends or explicitly follow random teleport links with probability 1.0 from dead-ends.

# The Google Matrix

- **PageRank equation** [Brin-Page, '98]

$$r_j = \sum_{i \to j} \beta \frac{r_i}{d_i} + (1 - \beta) \frac{1}{N}$$

- **The Google Matrix *A*:**

  $[1/N]_{NxN}$…N by N matrix where all entries are 1/N

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

- **We have a recursive problem: $r = A \cdot r$**
  **And the Power method still works!**
- **What is $\beta$ ?**
  - In practice $\beta = 0.8, 0.9$ (make $5$ steps on avg., jump)

# Random Teleports (β = 0.8)



**M**

$$0.8 \begin{vmatrix} 1/2 & 1/2 & 0 \\ 1/2 & 0 & 0 \\ 0 & 1/2 & 1 \end{vmatrix}$$

**[1/N]$_{NxN}$**

$$+ 0.2 \begin{vmatrix} 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \end{vmatrix}$$

|   | | | |
|---|---|---|---|
| y | 7/15 | 7/15 | 1/15 |
| a | 7/15 | 1/15 | 1/15 |
| m | 1/15 | 7/15 | 13/15 |

**A**

| y |   | 1/3 | 0.33 | 0.24 | 0.26 |     | 7/33 |
|---|---|-----|------|------|------|-----|------|
| a | = | 1/3 | 0.20 | 0.20 | 0.18 | . . . | 5/33 |
| m |   | 1/3 | 0.46 | 0.52 | 0.56 |     | 21/33 |

## Example Problem

Suppose we compute PageRank with a β of 0.7, and we introduce the additional constraint that the sum of the PageRanks of the three pages must be 3, to handle the problem that otherwise any multiple of a solution will also be a solution. Compute the PageRanks $a$, $b$, and $c$ of the three pages A, B, and C, respectively. Then, identify from the list below, the true statement.

a. $a + b = 1.025$
b. $a + b = 0.705$
c. $a + c = 2.035$
d. $a + b = 0.55$

## How do we actually compute the PageRank?

## Computing Page Rank

- **Key step is matrix-vector multiplication**
  - $r^{new} = A \cdot r^{old}$
- 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!

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

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

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

## Matrix Formulation

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

# Rearranging the Equation

- $r = A \cdot r$, where $A_{ji} = \beta M_{ji} + \frac{1-\beta}{N}$
- $r_j = \sum_{i=1}^{N} A_{ji} \cdot r_i$
- $r_j = \sum_{i=1}^{N} \left[ \beta M_{ji} + \frac{1-\beta}{N} \right] \cdot r_i$

  $= \sum_{i=1}^{N} \beta M_{ji} \cdot r_i + \frac{1-\beta}{N} \sum_{i=1}^{N} r_i$

  $= \sum_{i=1}^{N} \beta M_{ji} \cdot r_i + \frac{1-\beta}{N}$      since $\sum r_i = 1$

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

**Note:** Here we assumed **M** has no dead-ends

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

# Sparse Matrix Formulation

- We just rearranged the **PageRank equation**

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

  - where **[(1-β)/N]_N** is a vector with all **N** entries **(1-β)/N**

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

# PageRank: The Complete Algorithm

- **Input: Graph $G$ and parameter $\beta$**
  - Directed graph $G$ (can have **spider traps** and **dead ends**)
  - Parameter $\beta$
- **Output: PageRank vector $r^{new}$**

  - **Set:** $r_j^{old} = \frac{1}{N}$
  - **repeat until convergence:** $\sum_j \left| r_j^{new} - r_j^{old} \right| > \varepsilon$

    - $\forall j:$ $r'_j^{new} = \sum_{i \to j} \beta \frac{r_i^{old}}{d_i}$

      $r'_j^{new} = 0$  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-β**. But since we have dead-ends the amount of leaked PageRank may be larger. We have to explicitly account for it by computing **S**.