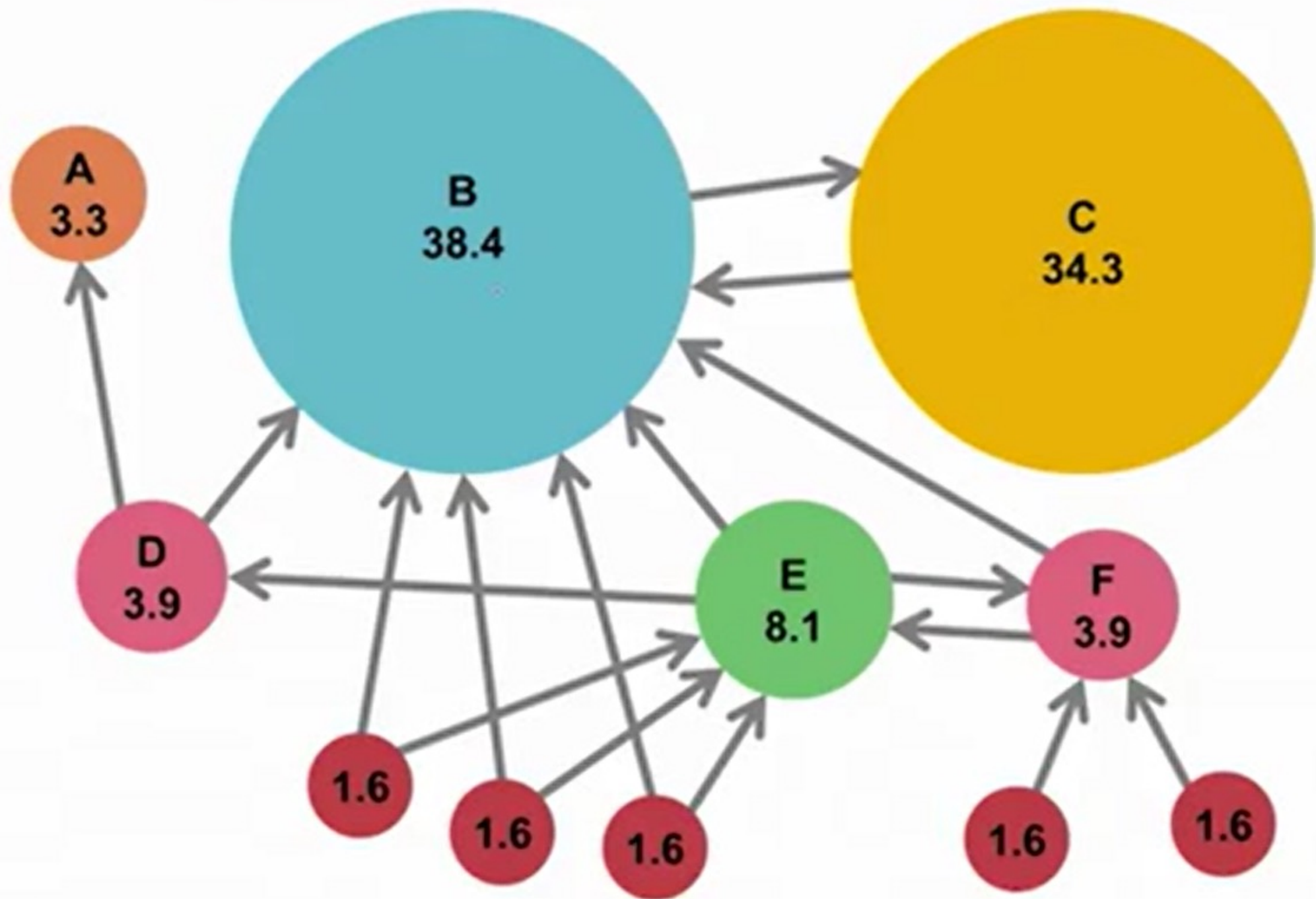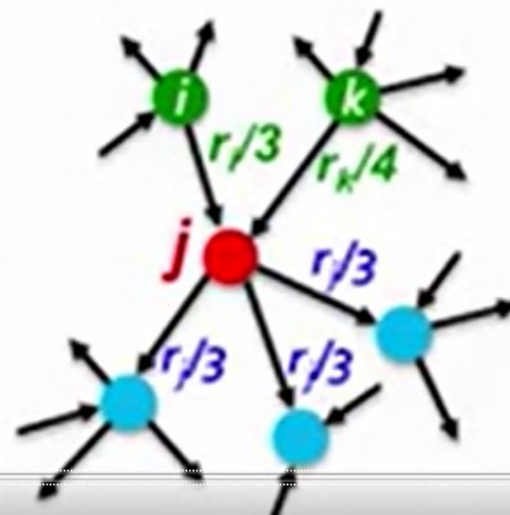# Page Rank

# Simple Recursive Formulation

- Each link's <u>vote</u> is proportional to the **importance** of its source page

- If page *j* with importance $r_j$ has **n** out-links, each link gets $r_j / n$ votes

- Page *j*'s own importance is the sum of the votes on its in-links
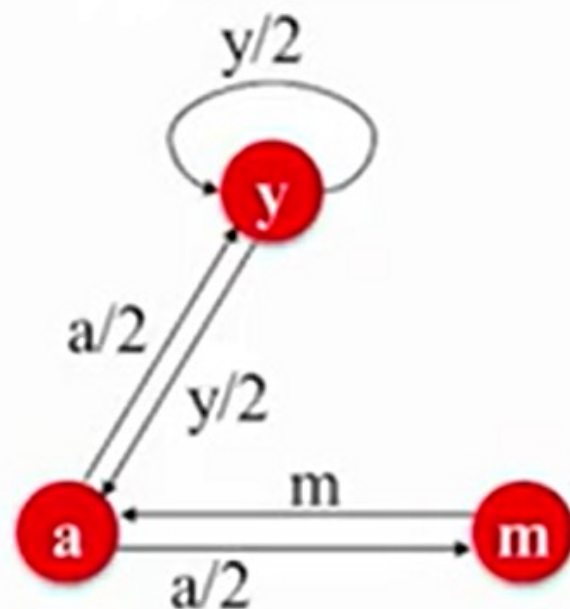
$$r_j = r_i/3 + r_k/4$$

# PageRank: The "Flow" Model

- A **"vote"** from an important page is worth more
- A page is important if it is pointed to by other important pages
- Define a "rank" $r_j$ for page $j$

$$r_j = \sum_{i \to j} \frac{r_i}{d_i}$$

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

# What is the flow equation?

Flow equations:
$$r_y = r_y/2 + r_a/2$$
$$r_a = r_y/2 + r_m$$
$$r_m = r_a/2$$

- **3 equations, 3 unknowns, no constants**
  - ⇨ No unique solution
  - All solutions equivalent modulo the scale factor
- **Additional constraint forces uniqueness:**

⇒ $r_y + r_a + r_m = 1$

⇒ Solution: $r_y = \dfrac{2}{5}$, $r_a = \dfrac{2}{5}$, $r_m = \dfrac{1}{5}$

- **Gaussian elimination method works for small examples, but we need a better method for large web-size graphs**
- **We need a new formulation!**

# PageRank: Matrix Formulation

- **Stochastic adjacency matrix $M$**
  - Let page $i$ has $d_i$ out-links
  - If $i \to j$, then $M_{ji} = \dfrac{1}{d_i}$ else $M_{ji} = 0$
    - $M$ is a **column stochastic matrix**
      - Columns sum to **1**

- **Rank vector $r$:** vector with an entry per page
  - $r_i$ is the importance score of page $i$
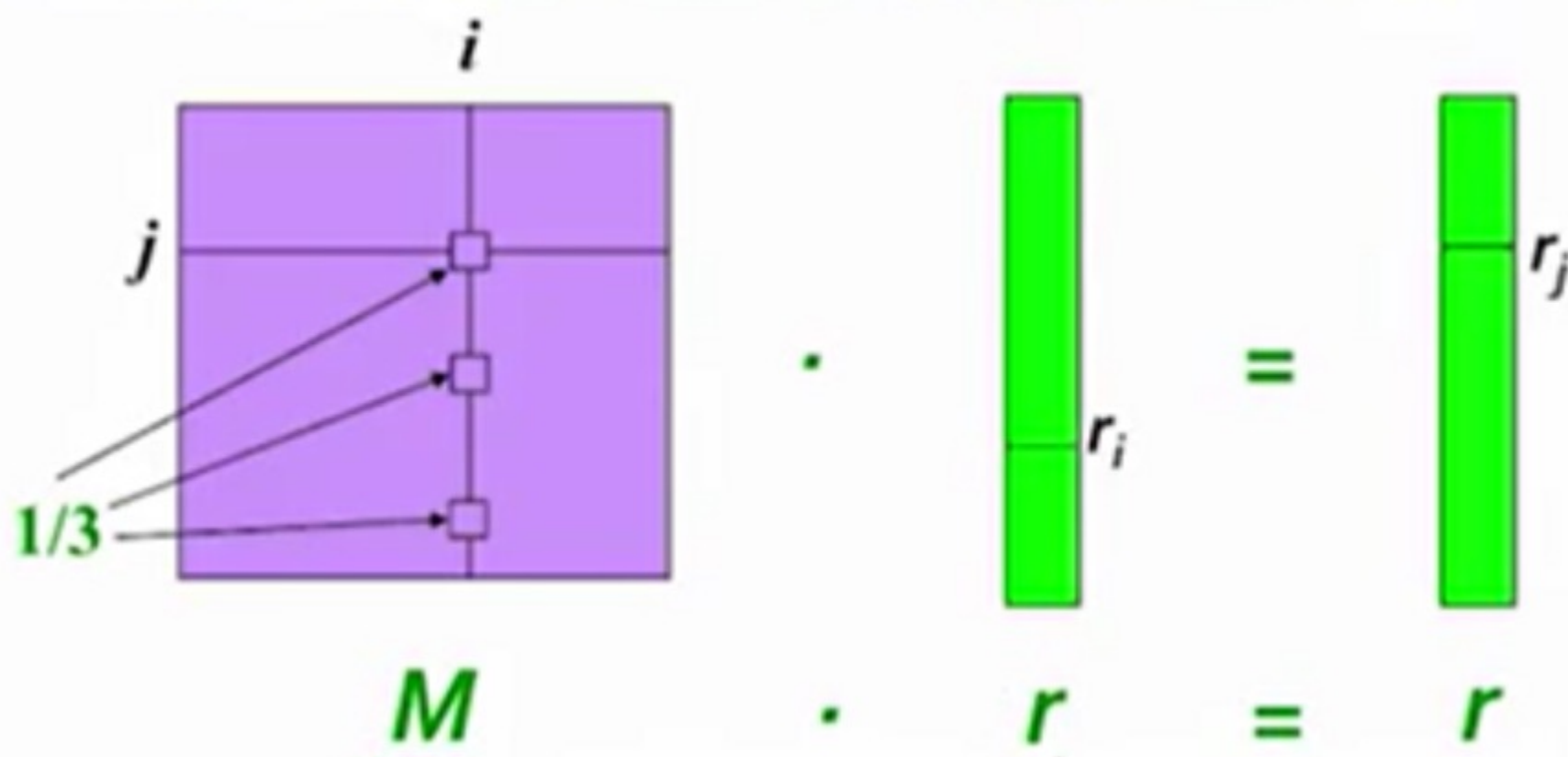  - $\to \sum_i r_i = 1$
- **The flow equations can be written**

$$r = M \cdot r$$

$$r_j = \sum_{i \to j} \frac{r_i}{d_i}$$

- **Remember the flow equation:** $r_j = \sum_{i \to j} \dfrac{r_i}{d_i}$
- **Flow equation in the matrix form**

$$M \cdot r = r$$

- **Suppose page $i$ links to 3 pages, including $j$**



$$M \qquad \cdot \qquad r \qquad = \qquad r$$

# Eigenvector Formulation
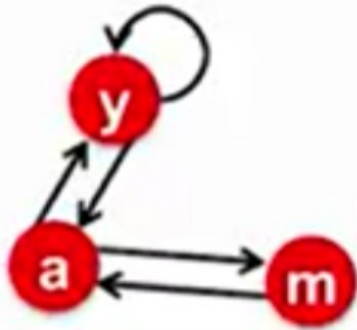
- **The flow equations can be written**

$$r = M \cdot r$$

  - So the **rank vector** *r* is an **eigenvector** of the stochastic web matrix **M**
    - In fact, its first or principal eigenvector, with corresponding eigenvalue *1*
      - Largest eigenvalue of **M** is **1** since **M** is column stochastic
        - **Why?** *We know* **r** *is a stochastic vector and each column of* M *sums to one, so* Mr ≤ 1

- **We can now efficiently solve for *r*!**
  **The method is called Power iteration**

# Flow equations and Matrix

$$M = \begin{array}{c|ccc} & y & a & m \\ \hline y & \tfrac{1}{2} & \tfrac{1}{2} & 0 \\ a & \tfrac{1}{2} & 0 & 1 \\ m & 0 & \tfrac{1}{2} & 0 \end{array}$$

$$r = M \cdot r$$

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

$$\begin{bmatrix} y \\ a \\ m \end{bmatrix} = \begin{bmatrix} \tfrac{1}{2} & \tfrac{1}{2} & 0 \\ \tfrac{1}{2} & 0 & 1 \\ 0 & \tfrac{1}{2} & 0 \end{bmatrix} \begin{bmatrix} y \\ a \\ m \end{bmatrix}$$

# Power Iteration Method

- **Given a web graph with *N* nodes, where the nodes are pages and edges are hyperlinks**
- **Power iteration:** a simple iterative scheme
  - Suppose there are *N* web pages
  - Initialize: $\mathbf{r}^{(0)} = [1/N,....,1/N]^{\mathsf{T}}$
  - Iterate: $\mathbf{r}^{(t+1)} = \mathbf{M} \cdot \mathbf{r}^{(t)}$
  - Stop when $|\mathbf{r}^{(t+1)} - \mathbf{r}^{(t)}|_1 < \varepsilon$

$$r_j^{(t+1)} = \sum_{i \to j} \frac{r_i^{(t)}}{d_i}$$
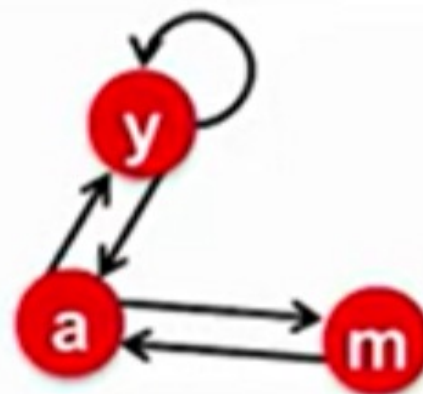
$d_i$ .... out-degree of node i

$|\mathbf{x}|_1 = \sum_{1 \le i \le N} |x|$ is the **L₁** norm
Can use any other vector norm, e.g., Euclidean

# PageRank: How to solve?

- **Power Iteration:**
  - Set $r_j = 1/N$
  - **1:** $r'_j = \sum_{i \to j} \frac{r_i}{d_i}$
  - **2:** $r = r'$
  - If not converged: goto **1**



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

$r_y = r_y/2 + r_a/2$

$r_a = r_y/2 + r_m$

$r_m = r_a/2$

- **Example:**

$$
\begin{pmatrix} r_y \\ r_a \\ r_m \end{pmatrix} =
\begin{matrix}
1/3 & 1/3 & 5/12 & 9/24 & & 6/15 \\
1/3 & 3/6 & 1/3 & 11/24 & \ldots & 6/15 \\
1/3 & 1/6 & 3/12 & 1/6 & & 3/15
\end{matrix}
$$

Iteration 0, 1, 2, ...