

Chapter IV:

Link Analysis

Information Retrieval & Data Mining
Universität des Saarlandes, Saarbrücken
Wintersemester 2013/14

Friendship Networks, Citation Networks, ...

- **Link analysis** studies the **relationships** (e.g., friendship, citation) between **objects** (e.g., people, publications) to find out about their **characteristics** (e.g., popularity, impact)
- **Social Network Analysis** (e.g., on a friendship network)
 - **Closeness centrality** of a person v is the **fraction of shortest paths** between any two persons (u, w) that pass through v
- **Bibliometrics** (e.g., on a citation network)
 - **Co-citation** measures how many papers cite both u and v
 - **Co-reference** measures how many common papers both u and v refer to

..., and the Web?

- **World Wide Web** can be seen as **directed graph** $G(V, E)$
 - **web pages** correspond to **vertices** (or, nodes) V
 - **hyperlinks** between them correspond to **edges** E
- Link analysis on the **Web graph** can give us clues about
 - which web pages are **important** and should thus be ranked higher
 - which pairs of web pages are **similar to each other**
 - which web pages are probably **spam** and should be ignored
 - ...

Chapter IV: Link Analysis

- IV.1 The World Wide Web as a Graph**
Degree Distributions, Diameter, Bow-Tie Structure
- IV.2 PageRank**
Random Surfer Model, Markov Chains
- IV.3 HITS**
Hyperlinked-Induced Topic Search
- IV.4 Topic-Specific and Personalized PageRank**
Biased Random Jumps, Linearity of PageRank
- IV.5 Online Link Analysis**
OPIC
- IV.6 Similarity Search**
SimRank, Random Walk with Restarts
- IV.7 Spam Detection**
Link Spam, TrustRank, SpamRank
- IV.8 Social Networks**
SocialPageRank, TunkRank

IV.1 The World Wide Web as a Graph

- 1. How Big is the Web?**
- 2. Degree Distributions**
- 3. Random-Graph Models**
- 4. Bow-Tie Structure**

Based on MRS Chapter 21

1. How Big is the Web?

- How big is the entire World Wide Web?
 - **quasi-infinite** when you consider all (dynamic) URLs (e.g., of calendars)
- **Indexed Web** is a more reasonable notion to look at
 - [Gulli and Signori '05] estimated it as 11.5 billions (10^9) in 2005
 - Google claimed to know about more than 1 trillion (10^{12}) URLs in 2008
 - WorldWideWebSize.com provides daily estimates obtained by extrapolating from the number of results returned by Google and Bing on the basis of Zipf's law (currently: 3.6 billion – 38 billion)

2. Degree Distributions

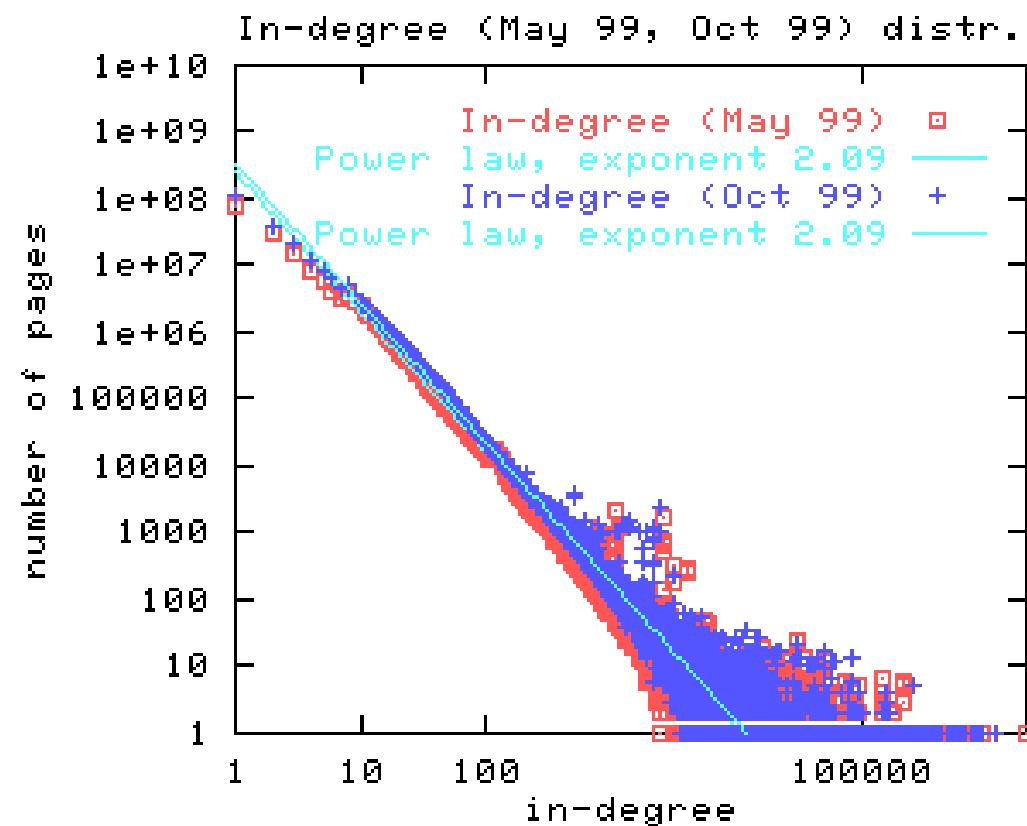
- What is the distribution of in-/out-degrees on the Web graph?
 - *in-degree*(v) of vertex v is the number of incoming edges (u, v)
 - *out-degree*(v) of vertex v is the number of outgoing edges (v, w)
- **Zipfian distribution** has probability mass function

$$f(k; s, N) = \frac{1/k^s}{\sum_{n=1}^N 1/n^s}$$

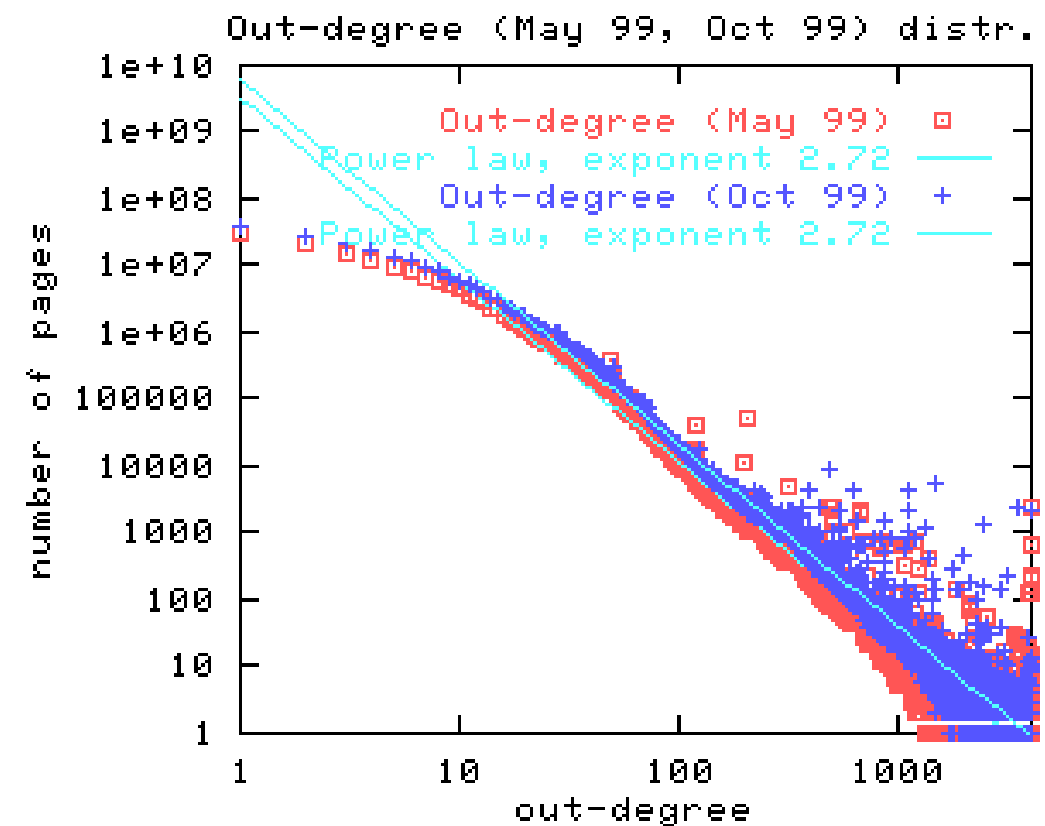
with rank k , parameter s , and total number of objects N

- provides good model of **many real-world phenomena**, e.g., word frequencies, city populations, corporation sizes, income rankings
- appear as **straight line** with slope $-s$ in **log-log-plot**

Degree Distributions



$$s = 2.10$$



$$s = 2.72$$

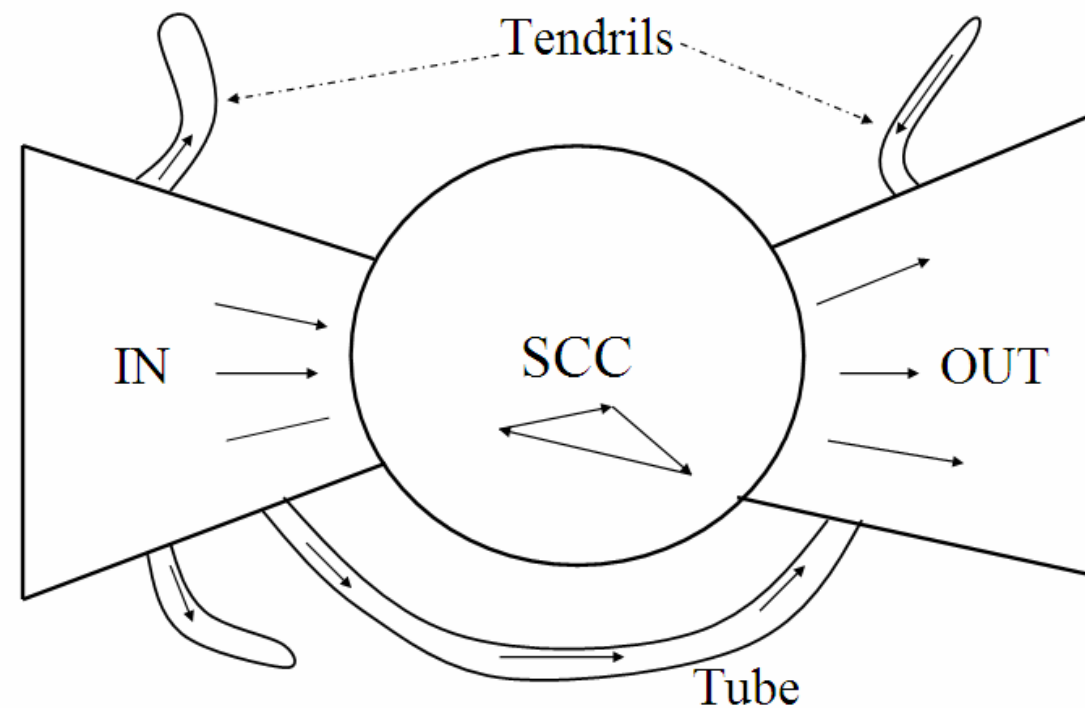
- Full details: [Broder et al. '00]

3. Random-Graph Models

- **Generative models** of undirected or undirected graphs
- **Erdős-Renyi Model** $G(n, p)$ generates a graph consisting of n vertices; each possible edge (u, w) exists with probability p
- **Barabási-Albert Model** generates a graph by successively adding vertices u with m edges; the edge (u, v) attaches to vertex v with probability proportional to $\deg(v)$
- **Preferential attachment** (“*the rich get richer*”) in the Barabási-Albert Model yields graphs with properties similar to Web graph
- Full details: [Barabási and Albert '99]

4. Bow-Tie Structure

- The Web graph looks a lot like a **bow tie** [Broder et al. '00]



- **Strongly Connected Component** (SCC) of web pages that are reachable from each other by following a few hyperlinks
- **IN** consisting of web pages from which SCC is reachable
- **OUT** consisting of web pages reachable from SCC

Additional Literature for IV.1

- **A.-L. Barabási and R. Albert:** *Emergence of Scaling in Random Networks*, Science 1999
- **A. Broder, R. Kumar, F. Maghoul, P. Raghavan, S. Rajagopalan, R. Stata, A. Tomkins, and J. L. Wiener:** *Graph Structure in the Web*, Computer Networks 33:309-320, 2000
- **A. Gulli and A. Signori:** *The Indexable Web is More than 11.5 Billion Pages*, WWW 2005
- **R. Meusel, O. Lehmberg, C. Bizer:** *Topology of the WDC Hyperlink Graph* <http://webdatacommons.org/hyperlinkgraph/topology.html>, 2013

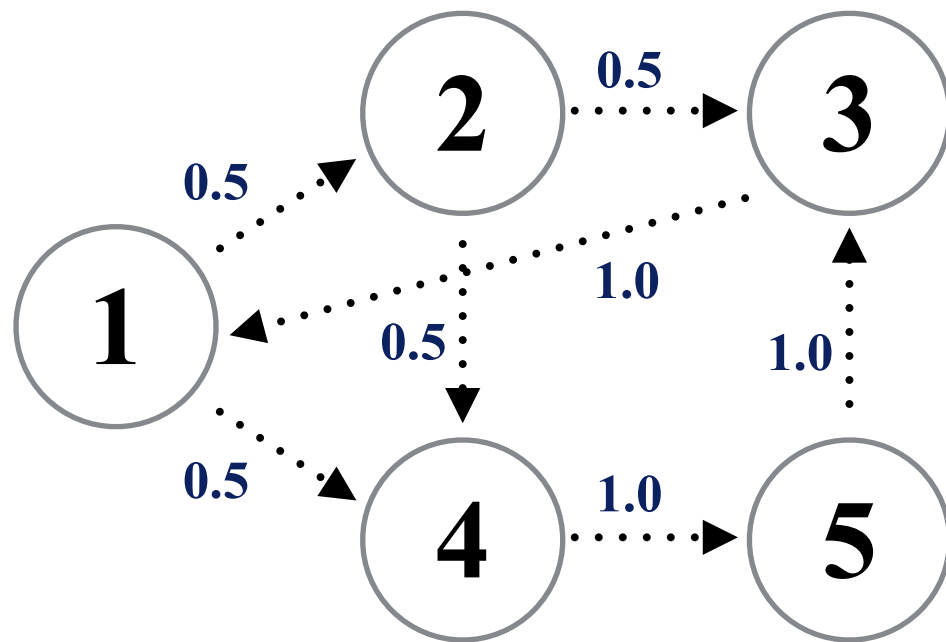
IV.2 PageRank

- **Hyperlinks** distinguish the Web from other document collections and can be interpreted as **endorsements** for the target web page
- **In-degree** as a measure of the **importance/authority/popularity** of a web page v is **easy to manipulate** and does not consider the **importance of the source web pages**
- **PageRank** considers a web page v **important** if **many important** web pages link to it
- **Random surfer model**
 - follows a uniform random outgoing link with probability $(1-\epsilon)$
 - jumps to a uniform random web page with probability ϵ
- Intuition: Important web pages are the ones that are visited often



Larry Page & Sergey Brin

Markov Chains



$$\mathbf{P} = \begin{bmatrix} 0.0 & 0.5 & 0.0 & 0.5 & 0.0 \\ 0.0 & 0.0 & 0.5 & 0.5 & 0.0 \\ 1.0 & 0.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0 & 1.0 \\ 0.0 & 0.0 & 1.0 & 0.0 & 0.0 \end{bmatrix}$$

$$S = \{1, \dots, 5\}$$

Stochastic Processes & Markov Chains

- **Discrete stochastic process** is a family of random variables

$$\{X_t \mid t \in T\}$$

with $T = \{0, 1, 2, \dots\}$ as **discrete time domain**

- Stochastic process is a **Markov chain** if

$$\begin{aligned} &P[X_t = x \mid X_{t-1} = w, \dots, X_0 = a] \\ &= P[X_t = x \mid X_{t-1} = w] \end{aligned}$$

holds, i.e., it is **memoryless**

- Markov chain is **time-homogeneous** if for all times t

$$P[X_{t+1} = x \mid X_t = w] = P[X_t = x \mid X_{t-1} = w]$$

holds, i.e., **transition probabilities** do not depend on time

State Space & Transition Probability Matrix

- **State space** of a Markov chain $\{X_t \mid t \in T\}$ is the countable set S of all values that X_t can assume
 - $X_t : \Omega \rightarrow S$
 - Markov chain **is in state s at time t** if $X_t = s$
 - Markov chain $\{X_t \mid t \in T\}$ is **finite** if it has a finite state space
- If a Markov chain $\{X_t \mid t \in T\}$ is **finite** and **time-homogeneous**, its **transition probabilities** can be described as a matrix $\mathbf{P} = (p_{ij})$

$$p_{ij} = P[X_t = j \mid X_{t-1} = i]$$

- For $|S| = n$ the transition probability matrix \mathbf{P} is a **n -by- n right-stochastic matrix** (i.e., its rows sum up to 1)

$$\forall i : \sum_j p_{ij} = 1$$

Properties of Markov Chains

- State i is **reachable** from state j if there exists a $n \geq 0$ such that $(\mathbf{P}^n)_{ij} > 0$ (with $\mathbf{P}^n = \mathbf{P} \times \dots \times \mathbf{P}$ as n -th exponent of \mathbf{P})
- States i and j **communicate** if i is reachable from j and vice versa
- Markov chain is **irreducible** if all states $i, j \in S$ communicate
- Markov chain is **positive recurrent** if the **recurrence probability** is 1 and the **mean recurrence time** is finite for every state i

$$\sum_{k=1}^{\infty} P[X_k = i \wedge \forall 1 \leq j < k : X_j \neq i \mid X_0 = i] = 1$$

$$\sum_{k=1}^{\infty} k P[X_k = i \wedge \forall 1 \leq j < k : X_j \neq i \mid X_0 = i] < \infty$$

Properties of Markov Chains

- Markov chain is **aperiodic** if every state i has period 1 defined as

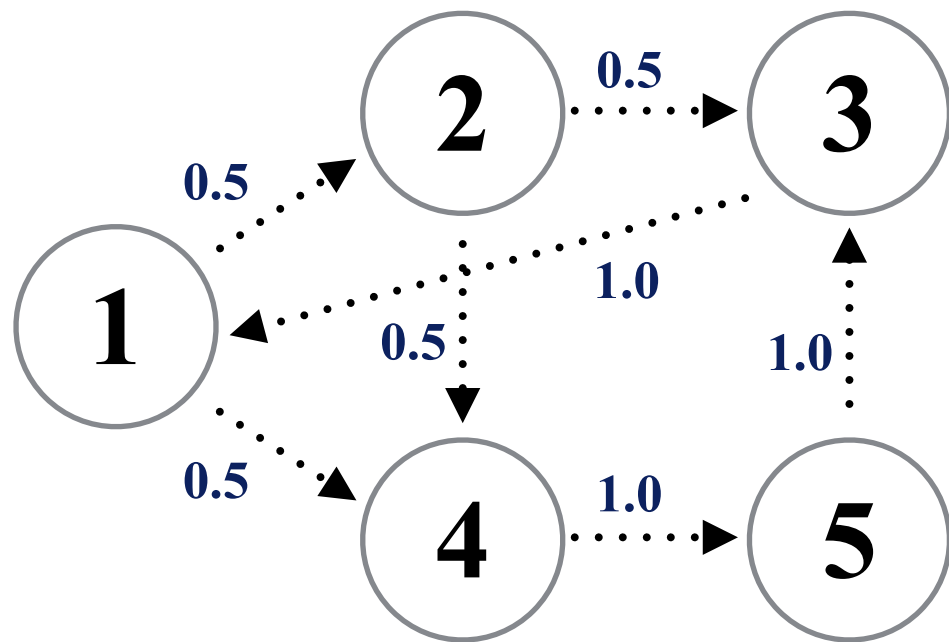
$$\gcd \{ k : P[X_k = i \wedge \forall 1 \leq j < k : X_j \neq i \mid X_0 = i] > 0 \}$$

- Markov chain is **ergodic** if it is time-homogeneous, irreducible, positive recurrent, and aperiodic
- The 1-by- n vector π is the **stationary state distribution** of the Markov chain described by \mathbf{P} if $\pi_i \geq 0$, $\sum \pi_i = 1$, and

$$\pi \mathbf{P} = \pi$$

- π_i is the limit probability that Markov chain is in state i
- $1/\pi_i$ reflects the average time until the Markov chain returns to state i
- Theorem: If a Markov chain is **finite** and **ergodic**, then there exists a **unique stationary state distribution** π

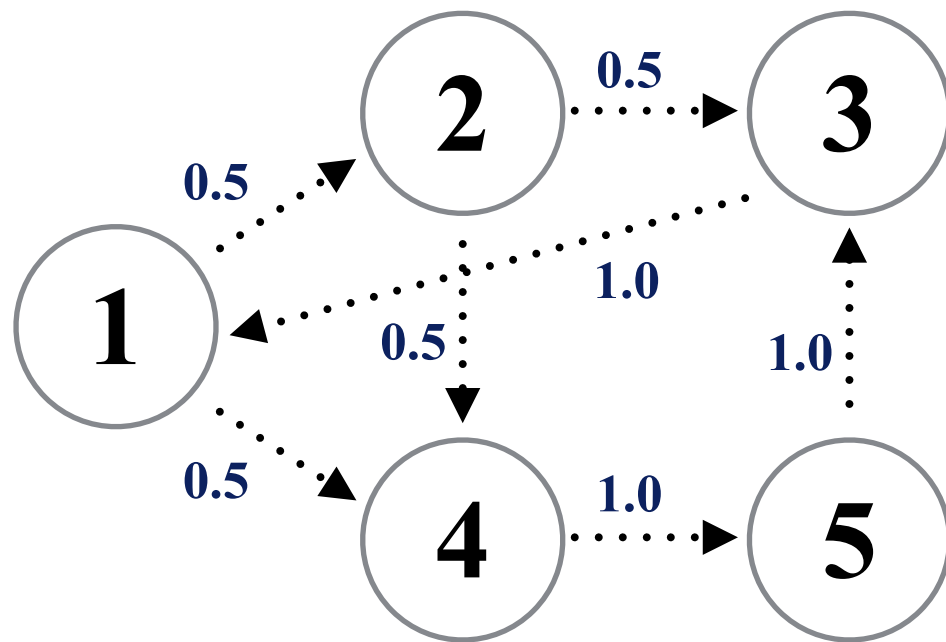
Markov Chain (Example Revisited)



$$S = \{1, \dots, 5\}$$

$$\mathbf{P} = \begin{bmatrix} 0.0 & 0.5 & 0.0 & 0.5 & 0.0 \\ 0.0 & 0.0 & 0.5 & 0.5 & 0.0 \\ 1.0 & 0.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0 & 1.0 \\ 0.0 & 0.0 & 1.0 & 0.0 & 0.0 \end{bmatrix}$$
$$\pi^0 = [1.0 \quad 0.0 \quad 0.0 \quad 0.0 \quad 0.0]$$

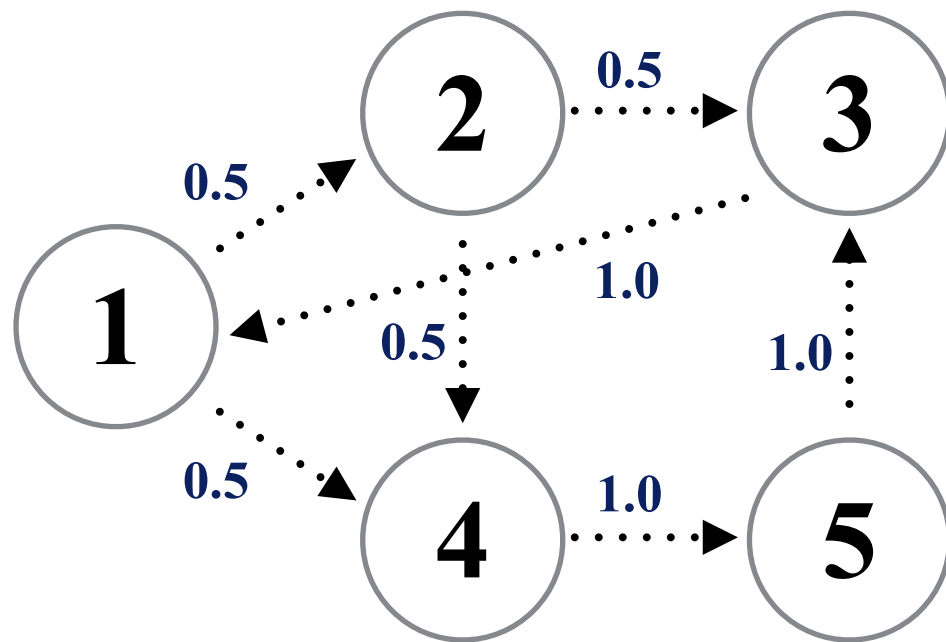
Markov Chain (Example Revisited)



$$\mathbf{P} = \begin{bmatrix} 0.0 & 0.5 & 0.0 & 0.5 & 0.0 \\ 0.0 & 0.0 & 0.5 & 0.5 & 0.0 \\ 1.0 & 0.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0 & 1.0 \\ 0.0 & 0.0 & 1.0 & 0.0 & 0.0 \end{bmatrix}$$
$$\pi^0 = [1.0 \quad 0.0 \quad 0.0 \quad 0.0 \quad 0.0]$$
$$\pi^1 = [0.0 \quad 0.5 \quad 0.0 \quad 0.5 \quad 0.0]$$

$$S = \{1, \dots, 5\}$$

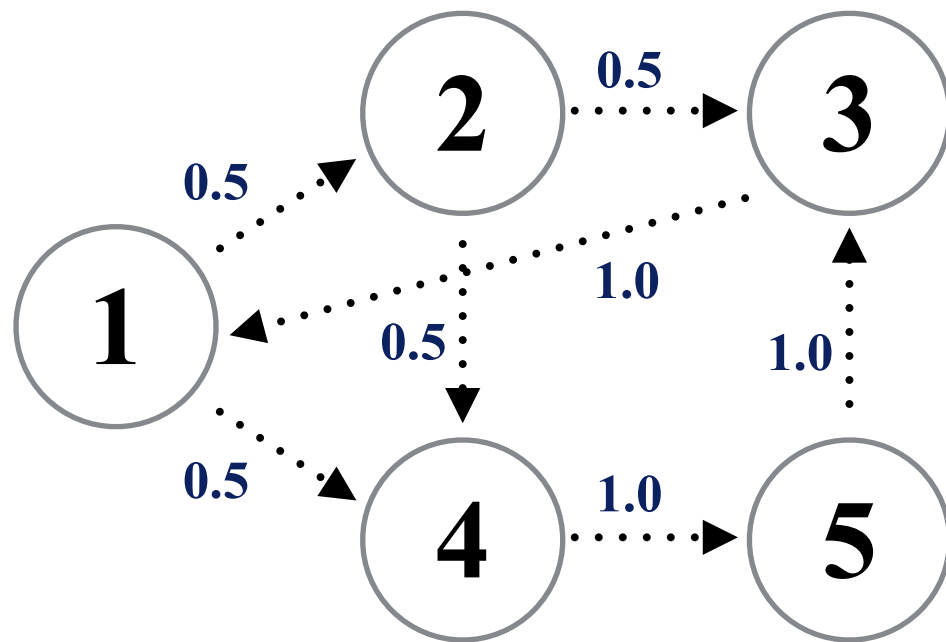
Markov Chain (Example Revisited)



$$S = \{1, \dots, 5\}$$

$$\mathbf{P} = \begin{bmatrix} 0.0 & 0.5 & 0.0 & 0.5 & 0.0 \\ 0.0 & 0.0 & 0.5 & 0.5 & 0.0 \\ 1.0 & 0.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0 & 1.0 \\ 0.0 & 0.0 & 1.0 & 0.0 & 0.0 \end{bmatrix}$$
$$\pi^0 = [1.0 \quad 0.0 \quad 0.0 \quad 0.0 \quad 0.0]$$
$$\pi^1 = [0.0 \quad 0.5 \quad 0.0 \quad 0.5 \quad 0.0]$$
$$\pi^2 = [0.0 \quad 0.0 \quad 0.25 \quad 0.25 \quad 0.5]$$

Markov Chain (Example Revisited)



$$S = \{1, \dots, 5\}$$

$$P = \begin{bmatrix} 0.0 & 0.5 & 0.0 & 0.5 & 0.0 \\ 0.0 & 0.0 & 0.5 & 0.5 & 0.0 \\ 1.0 & 0.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0 & 1.0 \\ 0.0 & 0.0 & 1.0 & 0.0 & 0.0 \end{bmatrix}$$

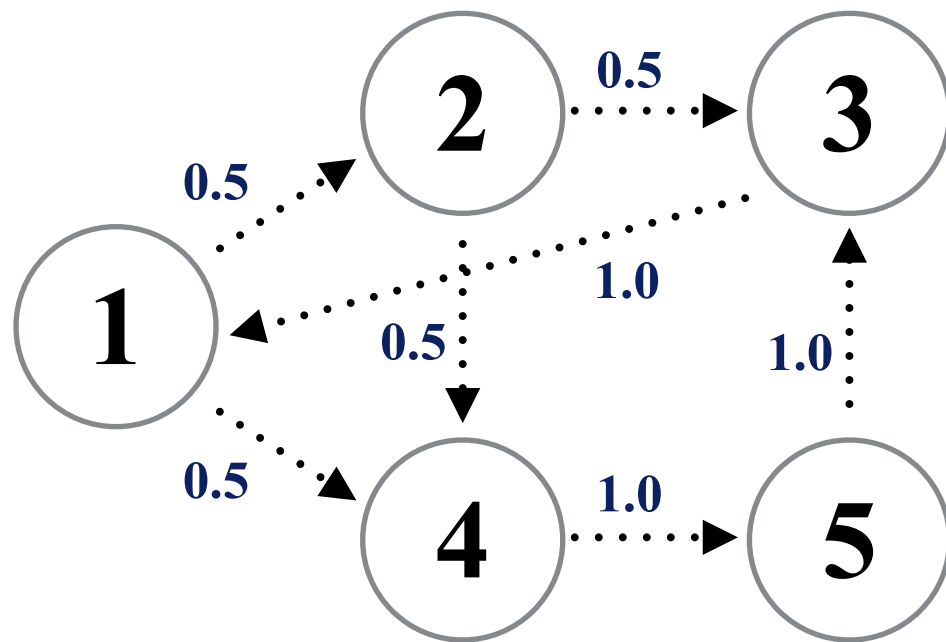
$$\pi^0 = [1.0 \quad 0.0 \quad 0.0 \quad 0.0 \quad 0.0]$$

$$\pi^1 = [0.0 \quad 0.5 \quad 0.0 \quad 0.5 \quad 0.0]$$

$$\pi^2 = [0.0 \quad 0.0 \quad 0.25 \quad 0.25 \quad 0.5]$$

$$\pi^3 = [0.25 \quad 0.0 \quad 0.5 \quad 0.0 \quad 0.25]$$

Markov Chain (Example Revisited)



$$S = \{1, \dots, 5\}$$

$$P = \begin{bmatrix} 0.0 & 0.5 & 0.0 & 0.5 & 0.0 \\ 0.0 & 0.0 & 0.5 & 0.5 & 0.0 \\ 1.0 & 0.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0 & 1.0 \\ 0.0 & 0.0 & 1.0 & 0.0 & 0.0 \end{bmatrix}$$

$$\pi^0 = [1.0 \quad 0.0 \quad 0.0 \quad 0.0 \quad 0.0]$$

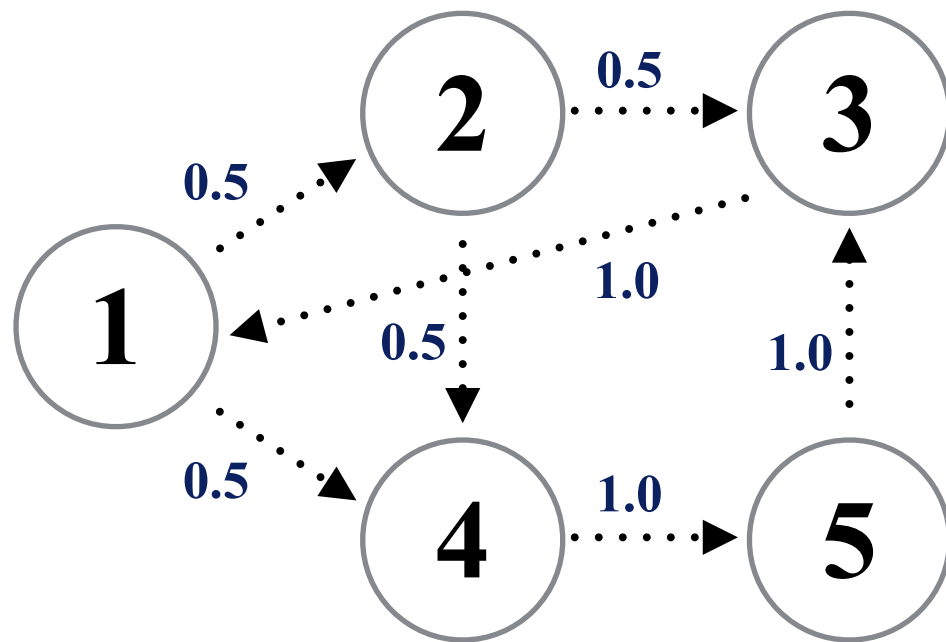
$$\pi^1 = [0.0 \quad 0.5 \quad 0.0 \quad 0.5 \quad 0.0]$$

$$\pi^2 = [0.0 \quad 0.0 \quad 0.25 \quad 0.25 \quad 0.5]$$

$$\pi^3 = [0.25 \quad 0.0 \quad 0.5 \quad 0.0 \quad 0.25]$$

$$\pi^4 = [0.5 \quad 0.125 \quad 0.25 \quad 0.125 \quad 0]$$

Markov Chain (Example Revisited)



$$S = \{1, \dots, 5\}$$

$$P = \begin{bmatrix} 0.0 & 0.5 & 0.0 & 0.5 & 0.0 \\ 0.0 & 0.0 & 0.5 & 0.5 & 0.0 \\ 1.0 & 0.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0 & 1.0 \\ 0.0 & 0.0 & 1.0 & 0.0 & 0.0 \end{bmatrix}$$

$$\pi^0 = [1.0 \quad 0.0 \quad 0.0 \quad 0.0 \quad 0.0]$$

$$\pi^1 = [0.0 \quad 0.5 \quad 0.0 \quad 0.5 \quad 0.0]$$

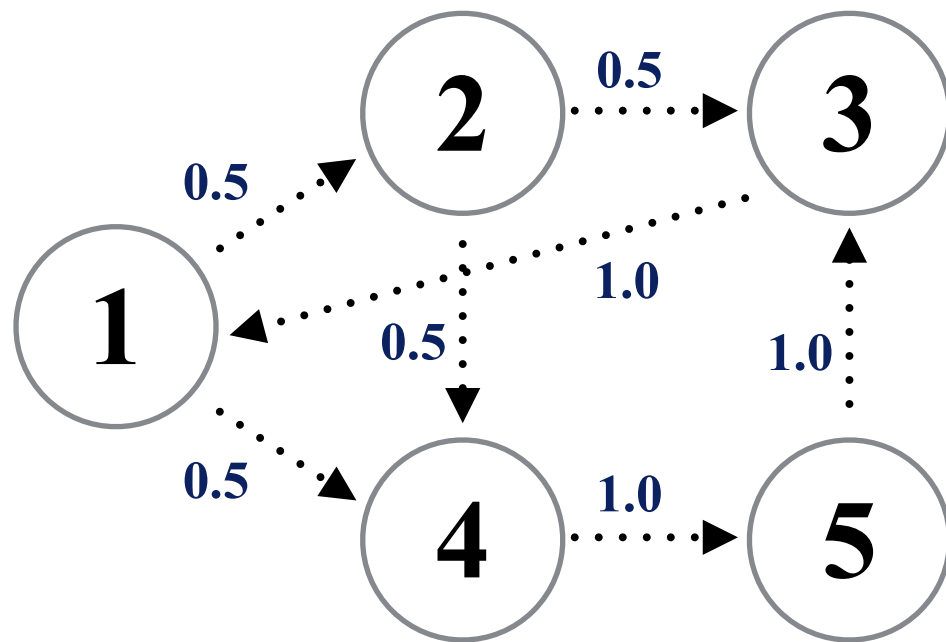
$$\pi^2 = [0.0 \quad 0.0 \quad 0.25 \quad 0.25 \quad 0.5]$$

$$\pi^3 = [0.25 \quad 0.0 \quad 0.5 \quad 0.0 \quad 0.25]$$

$$\pi^4 = [0.5 \quad 0.125 \quad 0.25 \quad 0.125 \quad 0]$$

$$\pi^5 = [0.25 \quad 0.25 \quad 0.0625 \quad 0.3125 \quad 0.125]$$

Markov Chain (Example Revisited)



$$S = \{1, \dots, 5\}$$

$$P = \begin{bmatrix} 0.0 & 0.5 & 0.0 & 0.5 & 0.0 \\ 0.0 & 0.0 & 0.5 & 0.5 & 0.0 \\ 1.0 & 0.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0 & 1.0 \\ 0.0 & 0.0 & 1.0 & 0.0 & 0.0 \end{bmatrix}$$

$$\pi^0 = [1.0 \quad 0.0 \quad 0.0 \quad 0.0 \quad 0.0]$$

$$\pi^1 = [0.0 \quad 0.5 \quad 0.0 \quad 0.5 \quad 0.0]$$

$$\pi^2 = [0.0 \quad 0.0 \quad 0.25 \quad 0.25 \quad 0.5]$$

$$\pi^3 = [0.25 \quad 0.0 \quad 0.5 \quad 0.0 \quad 0.25]$$

$$\pi^4 = [0.5 \quad 0.125 \quad 0.25 \quad 0.125 \quad 0]$$

$$\pi^5 = [0.25 \quad 0.25 \quad 0.0625 \quad 0.3125 \quad 0.125]$$

⋮

$$\pi = [0.25 \quad 0.125 \quad 0.25 \quad 0.1875 \quad 0.1875]$$

Computing π (Method 1)

- Stationary state distribution is the limit distribution
- Idea: Compute k -step state probabilities π^k until they converge
- **Power (iteration) method**
 - select arbitrary initial state probability distribution π^0
 - compute $\pi^k = \pi^{k-1} \mathbf{P}$ until they converge (e.g., $|\pi^k - \pi^{k-1}| < \varepsilon$)
 - report last π^k as stationary state distribution π
- Power (iteration) method basically **simulates the Markov chain** and is the **method of choice in practice** when dealing with huge state spaces, exploiting that **matrix-vector multiplication** is **easy to parallelize**

Computing π (Method 2)

- Stationary state distribution π fulfills $\pi = \pi P$, which can be cast into a **system of linear equations**

$$P = \begin{bmatrix} 0.0 & 0.5 & 0.0 & 0.5 & 0.0 \\ 0.0 & 0.0 & 0.5 & 0.5 & 0.0 \\ 1.0 & 0.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0 & 1.0 \\ 0.0 & 0.0 & 1.0 & 0.0 & 0.0 \end{bmatrix}$$

$$\pi_1 = 0.0 \times \pi_1 + 0.0 \times \pi_2 + 1.0 \times \pi_3 + 0.0 \times \pi_4 + 0.0 \times \pi_5$$

$$\pi_2 = 0.5 \times \pi_1 + 0.0 \times \pi_2 + 0.0 \times \pi_3 + 0.0 \times \pi_4 + 0.0 \times \pi_5$$

$$\pi_3 = 0.0 \times \pi_1 + 0.5 \times \pi_2 + 0.0 \times \pi_3 + 0.0 \times \pi_4 + 1.0 \times \pi_5$$

$$\pi_4 = 0.5 \times \pi_1 + 0.5 \times \pi_2 + 0.0 \times \pi_3 + 0.0 \times \pi_4 + 0.0 \times \pi_5$$

$$\pi_5 = 0.0 \times \pi_1 + 0.0 \times \pi_2 + 0.0 \times \pi_3 + 1.0 \times \pi_4 + 0.0 \times \pi_5$$

$$1 = 1.0 \times \pi_1 + 1.0 \times \pi_2 + 1.0 \times \pi_3 + 1.0 \times \pi_4 + 1.0 \times \pi_5$$

- Solutions can be found, e.g., using **Gauss elimination**

Computing π (Method 3)

- Stationary state probability distribution π is the left **eigenvector** of the **transition probability matrix** \mathbf{P} for the **eigenvalue** $\lambda = 1$

$$\pi \mathbf{P} = \lambda \pi$$

- Can be computed using the **characteristic polynomial**

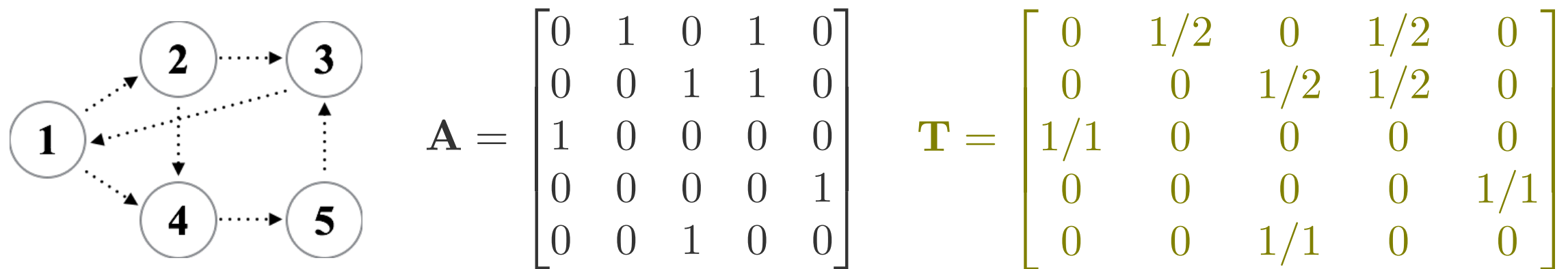
$$(\mathbf{P} - \lambda \mathbf{I}) \pi = 0$$

PageRank as a Markov Chain

- **Random surfer model**

- follows a uniform random outgoing link with probability $(1-\varepsilon)$
- jumps to a uniform random web page with probability ε
- Let **A** be the **adjacency matrix** of the Web graph, matrix **T** captures following of a uniform random outgoing link

$$\mathbf{T}_{ij} = \begin{cases} 1/out(i) & : (i, j) \in E \\ 0 & : \text{otherwise} \end{cases}$$

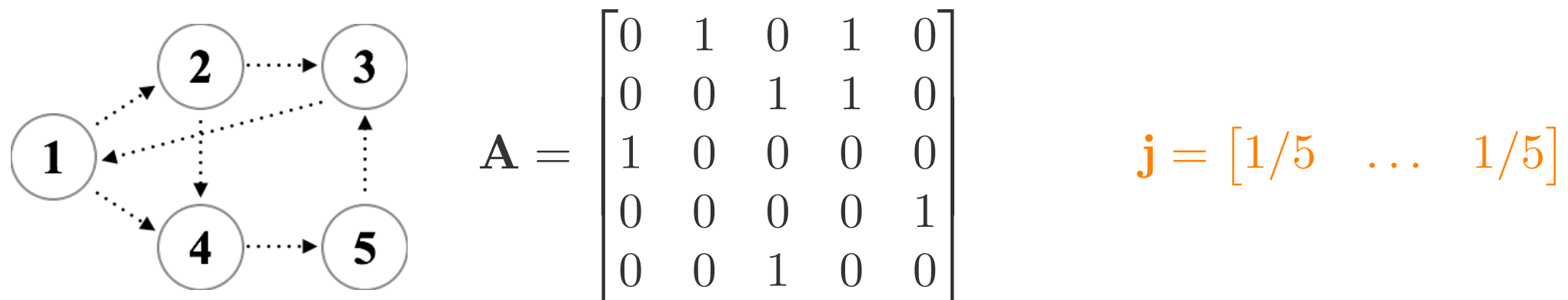


PageRank as a Markov Chain

- **Random surfer model**

- follows a uniform random outgoing link with probability $(1-\epsilon)$
- jumps to a uniform random web page with probability ϵ
- Vector \mathbf{j} captures jumping to a uniform random web page

$$\mathbf{j}_i = 1/|V|$$



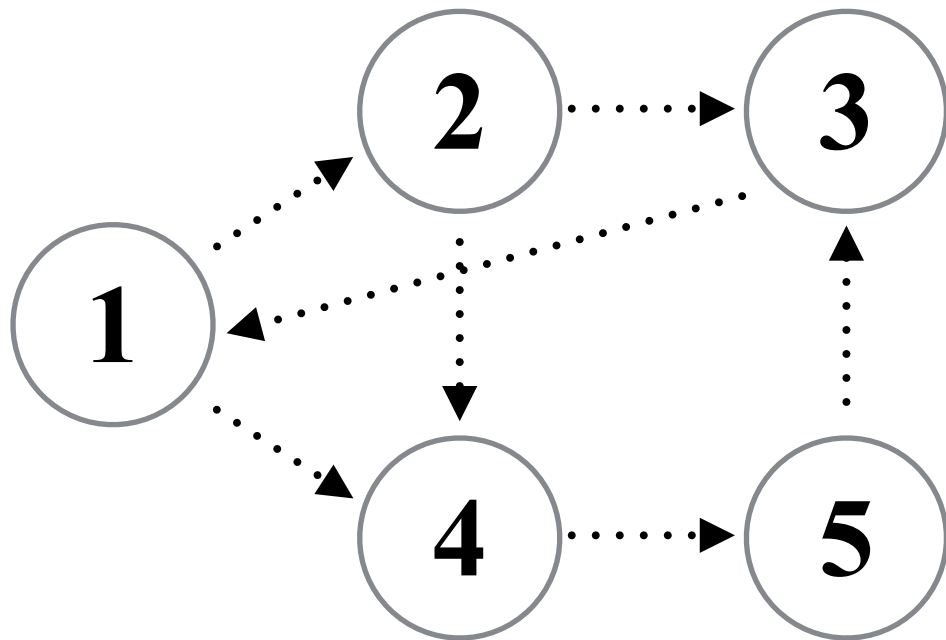
- **Transition probability matrix** of Markov chain then obtained as

$$\mathbf{P} = (1 - \epsilon) \mathbf{T} + \epsilon \begin{bmatrix} 1 & \dots & 1 \end{bmatrix}^T \mathbf{j}$$

PageRank as a Markov Chain

- With $\varepsilon = 0.15$ we obtain

$$\mathbf{P} = \begin{bmatrix} 0.030 & 0.455 & 0.030 & 0.455 & 0.030 \\ 0.030 & 0.030 & 0.455 & 0.455 & 0.030 \\ 0.880 & 0.030 & 0.030 & 0.030 & 0.030 \\ 0.030 & 0.030 & 0.030 & 0.030 & 0.880 \\ 0.030 & 0.030 & 0.880 & 0.030 & 0.030 \end{bmatrix}$$



$$\pi = [0.24079 \quad 0.13234 \quad 0.24799 \quad 0.18858 \quad 0.19029]$$

PageRank as a Markov Chain

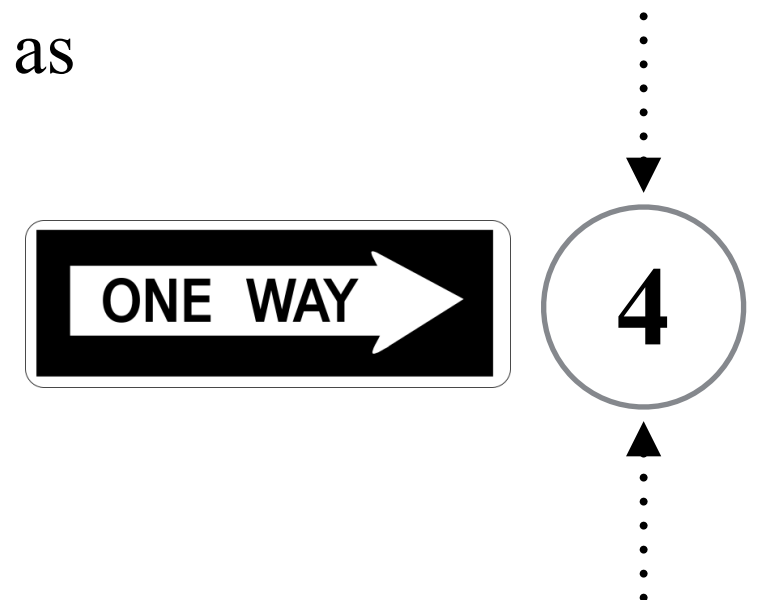
- **Transition probability matrix** of Markov chain then obtained as

$$\mathbf{P} = (1 - \epsilon) \mathbf{T} + \epsilon \begin{bmatrix} 1 & \dots & 1 \end{bmatrix}^T \mathbf{j}$$

$$\pi_i = (1 - \epsilon) \sum_{(j,i) \in E} \frac{\pi_j}{out(j)} + \frac{\epsilon}{|V|}$$

- We need to deal with **dangling nodes** (having out-degree zero)
- **Re-normalize** π^k such that $|\pi^k| = 1$ after every iteration of power method
- Make \mathbf{P} truly right stochastic by defining matrix \mathbf{T} as

$$\mathbf{T}_{ij} = \begin{cases} 1/out(i) & : (i, j) \in E \\ 1/|V| & : out(i) = 0 \\ 0 & : \text{otherwise} \end{cases}$$



PageRank as a Markov Chain (Is It Ergodic?)

- Markov chain defined by transition probability matrix \mathbf{T} is
 - **finite** (only finite number of web pages)
 - **time-homogeneous** (by design)
 - **irreducible** (random surfer can jump from every state i to every state j)
 - **positive recurrent** (random surfer can “jump up” on state i)
 - **aperiodic** (period of every state is 1 because of “jump up” on state i)

...it is thus **ergodic** and unique stationary state probabilities π exist

- **Random jump** is **essential** to make the Markov chain ergodic

PageRank & Queries

- **Random jump probability** typically set as $\varepsilon = 0.15$
(i.e., random surfer follows on average about seven links in a row)
- PageRank determines a **static global ranking** of web pages, is **query-independent**, and **orthogonal to textual relevance**

- Combination of PageRank score and retrieval models, e.g., as

- **linear combination** of cosine similarity and PageRank score

$$\alpha \times \text{sim}(q, d) + (1 - \alpha) \times \text{pr}(d)$$

- **document prior** in a query-likelihood language model

$$P(q|d) \times P(d)$$

- together with many other features in **machine-learned ranking model**

Summary of IV.2

- **Markov chains**

as a kind of stochastic process useful to describe random walks

- **Stationary state distribution**

is guaranteed to exist if the Markov chain is finite and ergodic
can be computed using (i) power iteration (ii) solving a system of linear equations or (iii) determining an eigenvector of a matrix

- **PageRank**

as Google's initial secret of success

is based on a random surfer model

can be described as a finite and ergodic Markov chain

yields a static query-independent importance score

Additional Literature for IV.2

- **S. Brin and L. Page:** *The anatomy of a large-scale hypertextual Web search engine*, Computer Networks 30:107-117, 1998
- **M. Bianchini, M. Gori, and F. Scarselli:** *Inside PageRank*, ACM TOIT 5(1):92-128, 2005
- **M. Franceschet:** *PageRank: Standing on the Shoulders of Giants*, CACM 54(6):92-101, 2011
- **A. N. Meyer and C. D. Meyer:** *Survey: Deeper Inside PageRank*, Internet Mathematics 1(3):335-380, 2003
- **L. Page, S. Brin, R. Motwani, and T. Winograd:** *The PageRank Citation Ranking: Bringing Order to the Web*, Technical Report, Stanford University, 1999