

# Machine Learning on Graphs

## MDI343

## PageRank

Thomas Bonald

2020 – 2021



# Motivation

How to identify the most “important” nodes in a graph, either **globally** or **relatively** to some other nodes?

Useful for:

- ▶ information retrieval
- ▶ content recommendation
- ▶ local clustering

We focus on **PageRank**, originally proposed by Google’s founders in 1999 to rank Web pages: popular pages are typically visited more frequently by a random Web surfer.

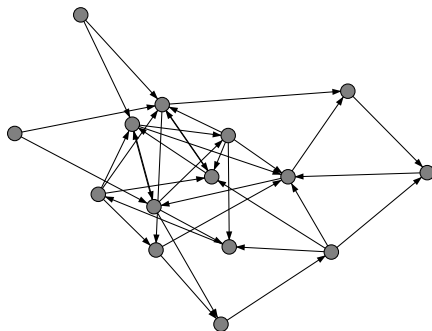
# Outline

1. Random walk
2. PageRank
3. Personalized PageRank

# Setting

Consider a directed graph  $G = (V, E)$ :

- ▶  $n$  nodes,  $m$  edges
- ▶  $A$ , adjacency matrix
- ▶  $d^+ = A\mathbf{1}$ ,  $d^- = A^T\mathbf{1}$ , vectors of out-degrees and in-degrees



# Random walk

In the **absence** of sinks ( $d^+ > 0$ ):

- ▶ A Markov chain  $X_0, X_1, X_2, \dots$  of transition matrix  $P = D^{-1}A$  with  $D = \text{diag}(d^+)$
- ▶ Probability distribution  $\pi_t$  at time  $t$  (row vector)
- ▶ Dynamics  $\pi_{t+1} = \pi_t P$

## Stationary distribution

If the graph is **strongly connected** and **aperiodic**,

$$\lim_{t \rightarrow +\infty} \pi_t = \pi \quad \text{with} \quad \pi = \pi P$$

# Computation

Stationary distribution

**Input:**

$P$ , transition matrix

$K$ , number of iterations

**Do:**

$\pi \leftarrow \frac{1}{n}(1, \dots, 1)$

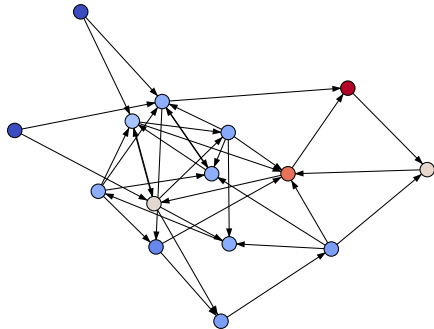
For  $t = 1, \dots, K$ ,  $\pi \leftarrow \pi P$

**Output:**

$\pi$ , (approximate) stationary distribution

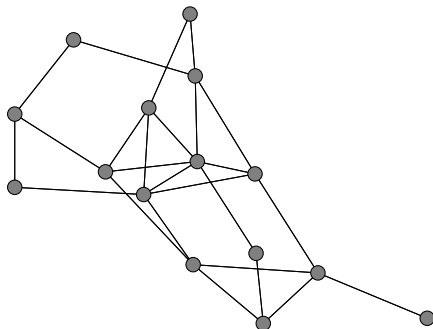
Complexity:  $O(Km)$  in time,  $O(n)$  in memory

## Example



# The case of undirected graphs

We have  $d = d^+ = d^-$



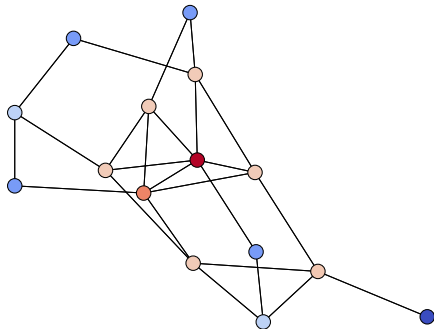
## Stationary distribution

If the graph is **connected**, the stationary distribution is proportional to the degrees:

$$\pi \propto d$$



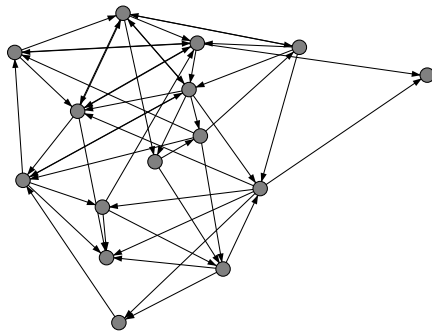
## Example



# Outline

1. Random walk
2. **PageRank**
3. Personalized PageRank

## Accounting for sinks



Random walk with forced restarts

$$P_{ij} = \begin{cases} \frac{A_{ij}}{d_i^+} & \text{if } d_i^+ > 0 \\ \frac{1}{n} & \text{otherwise} \end{cases}$$

# PageRank

Random walk with **restarts**:

- ▶ Fix  $\alpha \in (0, 1)$
- ▶ Walk with probability  $\alpha$ , restart (e.g., to a random node) with probability  $1 - \alpha$
- ▶ An irreducible Markov chain with transition matrix:

$$P^{(\alpha)} = \alpha P + (1 - \alpha) \frac{11^T}{n}$$

## PageRank

Unique solution to the equations:

$$\pi^{(\alpha)} = \alpha \pi^{(\alpha)} P + (1 - \alpha) \frac{1^T}{n}$$

# Computation

## PageRank

### **Input:**

$P$ , transition matrix (with forced restarts)

$\alpha$ , damping factor

$K$ , number of iterations

### **Do:**

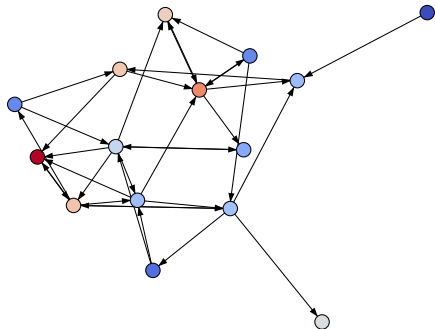
$$\pi \leftarrow \frac{1}{n}(1, \dots, 1)$$

$$\text{For } t = 1, \dots, K, \pi \leftarrow \alpha \pi P + (1 - \alpha) \frac{1}{n}(1, \dots, 1)$$

### **Output:**

$\pi$ , (approximate) PageRank vector

Example ( $\alpha = 0.85$ )



## Setting the damping factor

- ▶ The path length before restart (in the absence of sinks) has a **geometric distribution** with parameter  $1 - \alpha$

- ▶ Average path length:

$$\frac{\alpha}{1 - \alpha}$$

- ▶ For  $\alpha = 0.85$ , we get about 5.7, a typical distance between two nodes in real graphs (cf. the **six degrees of separation**).

# Expression of the PageRank vector

## Proposition

$$\pi^{(\alpha)} = (1 - \alpha) \sum_{t=0}^{+\infty} \alpha^t \pi_t$$

## Limiting cases

- ▶ **No restarts** ( $\alpha \rightarrow 1$ )

$$\pi^{(\alpha)} \rightarrow \pi = \lim_{t \rightarrow +\infty} \pi_t$$

- ▶ **Frequent restarts** ( $\alpha \rightarrow 0$ )

$$\pi^{(\alpha)} = (1 - \alpha)\pi_0 + \alpha\pi_1 + o(\alpha)$$

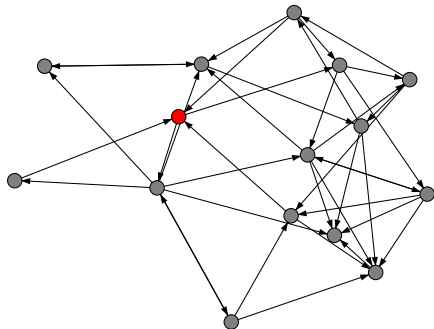
Ranking equivalent to neighbor sampling



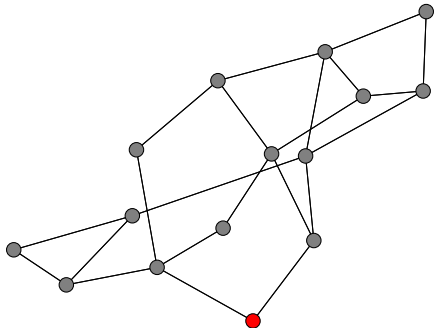
# Outline

1. Random walk
2. PageRank
3. **Personalized PageRank**

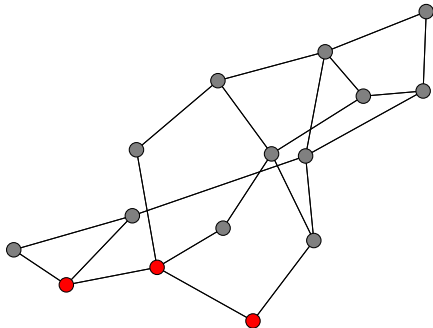
# Personalization



# Personalization



# Personalization



# Personalized PageRank

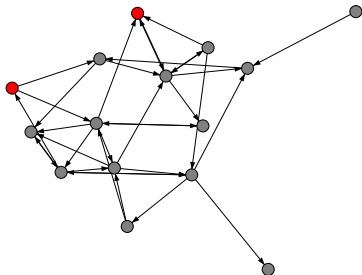
Let  $\mu$  be some distribution on  $S \subset V$  (e.g., uniform)

- Forced restarts:

$$P_{ij} = \begin{cases} \frac{A_{ij}}{d_i^+} & \text{if } d_i^+ > 0 \\ \mu_j & \text{otherwise} \end{cases}$$

- Random restarts:

$$P^{(\alpha)} = \alpha P + (1 - \alpha)1\mu$$



# Computation

## Personalized PageRank

### **Input:**

$P$ , transition matrix (with forced restarts)

$\mu$ , personalization row vector

$\alpha$ , damping factor

$K$ , number of iterations

### **Do:**

$\pi \leftarrow \mu$

For  $t = 1, \dots, K$ ,  $\pi \leftarrow \alpha \pi P + (1 - \alpha) \mu$

### **Output:**

$\pi$ , (approximate) PageRank vector

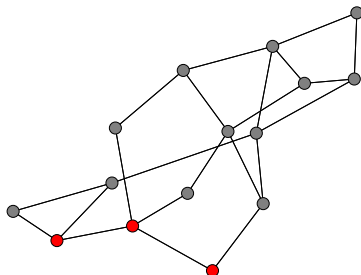
# Expression of the Personalized PageRank vector

## Proposition

In the absence of sinks,

$$\pi^{(\alpha)} = \sum_{s \in S} \mu_s \pi_s^{(\alpha)}$$

where  $\pi_s^{(\alpha)}$  is the Personalized PageRank vector associated with  $s$



# Summary

PageRank is a **key tool** for graph analysis:

- ▶ Useful to quantify the importance of nodes, possibly relatively to other nodes → **Personalized PageRank**
- ▶ **Fast** computation through matrix-vector multiplications

