# Exercise 10

Information Retrieval

# 13. Web Search Basics
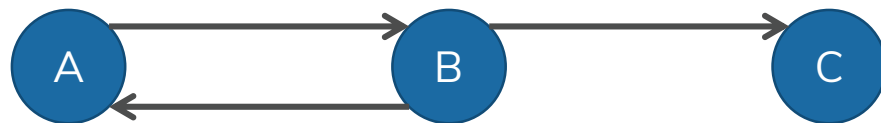# 14. PageRank and HITS

# Warm up

## *Exercise 13/14.1*

- Are the following statements true or false? Give reasons for your answer.

  a) The size of the web can easily be determined by crawling the web and counting the pages.

  b) In the context of web search, information needs are the only subclass of user needs. other subsets

  c) Shingling is a technique for detecting near duplicates of web pages.

  d) The average out-degree of all web pages is higher than their average in-degree.

  e) The PageRank algorithm ranks web pages by the number of occurrences of the query terms.

  f) The PageRank of a web page is query-dependent.

  g) The HITS-algorithm provides two scores per web page: an authority score and a hub score.

  h) The HITS-scores of a web page are query-dependent.

# PageRank: Transition Probability Matrices

## Exercise 13/14.2

- Consider the web graph shown below



- Write down the transition probability matrices for the random surfer's walk with teleporting
- Consider the following three values of the teleport probability

$$\alpha = 0 \qquad \alpha = \frac{1}{2} \qquad \alpha = 1$$

# PageRank and HITS

*Exercise 13/14.3*

- For the web graph shown on the right, compute the following for each of the three pages A, B, and C
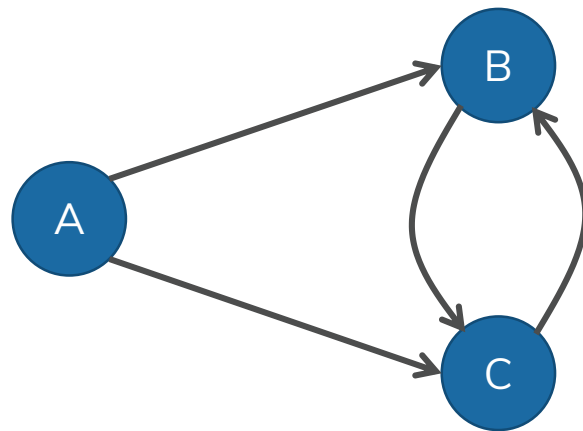
| PageRank | authority score | hub score |
|---|---|---|

- Also give the relative ordering of the three nodes for each of these scores, indicating any ties

*Hints*

1.
2.

- <u>PageRank</u>: Assume that at each step of the random walk, we teleport to a random page with probability 0.1, with a uniform distribution over which particular page we teleport to

- <u>Hubs/Authorities</u>: Normalize the hub (authority) scores so that the maximum hub (authority) score is 1

# Minimum PageRank

*Exercise 13/14.4*

a) Show that the PageRank of every page is at least $\alpha/n$

   where $\alpha$ is the teleport probability and $n$ is the total number of web pages

b) What does this imply regarding the difference in PageRank values (over the various pages) as $\alpha$ becomes close to 1? random i think