

# Markov Chains

## WEBSITE TRAFFIC PREDICTION

TANAY BIRADAR

OCTOBER 2, 2021

# WEBSITE TRAFFIC AND PAGERANK

- Search engines use popularity to rank pages
- Popularity can be quantified by links to a page
  - ▶ Works in theory, can be abused in practice
- Google used a Markov Model (PageRank) to rank popularity

# WEBSITE TRAFFIC AND PAGERANK

“PageRank works by counting the number and quality of links to a page to determine a rough estimate of how important the website is. The underlying assumption is that more important websites are likely to receive more links from other websites.” <sup>1</sup>

---

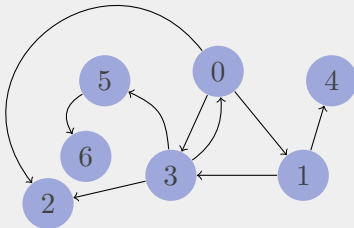
<sup>1</sup><https://en.wikipedia.org/wiki/PageRank>

# CENTRAL QUESTION

Which page is a user most likely to land on after starting on a given page?

# MODEL

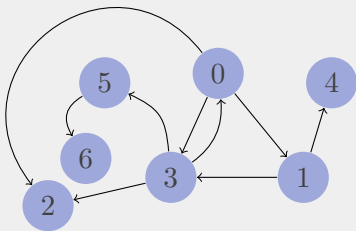
- Represent the internet as a directed graph
  - ▶ We're looking at a small slice of the web
- Edges are links, vertices are web pages
  - ▶ Assume equal probability of traversing every link such that  $\sum w_{out} = 1$ , where  $w$  is the edge weight
    - The probabilities coming out of every website must sum to 1



<sup>2</sup><https://en.wikipedia.org/wiki/PageRank>

The most complex part is by far the data collection

- Start at an arbitrary website, perform BFS to create an adjacency list that represents the internet
  - ▶ Keep track of visited nodes to avoid duplicate processing
- Stop after storing -- thousand links
  - ▶ I don't have the computing power of Google



# DATA (CONT'D)

Adjacency list  $A$  stores links between pages

If  $A_{ij} = 1$ , there is a link from page  $i$  to page  $j$

$$A = \begin{bmatrix} a_{00} & \dots & a_{0n} \\ \vdots & \ddots & \vdots \\ a_{n0} & \dots & a_{nn} \end{bmatrix}$$

Normalize the adjacency list to satisfy  $\forall i \sum w_{out} = \sum w_i = 1$

We now have a transition matrix!

Start at a given page