

L23
~~L22:~~ PageRank

Jeff M. Phillips

April 11, 2018

Final Report

At most 4 pages/student. Don't cram in too much!

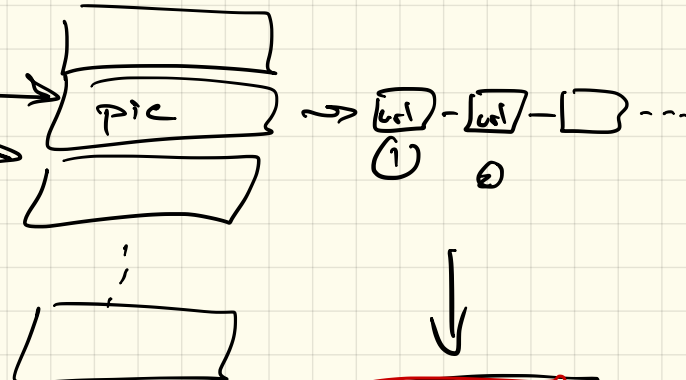
- ▶ Succinct title (and names)
- ▶ Problem definition and motivation.
- ▶ Explain your Data.
- ▶ **key idea**
- ▶ What did you do (which techniques, an implementation, a comparison, an extension)
- ▶ What did you learn? Artifacts (charts, plots, examples, math) and Intuition (in words, did it work?)

Page Rank ← key technology inside Google Search.

• Search Engine?

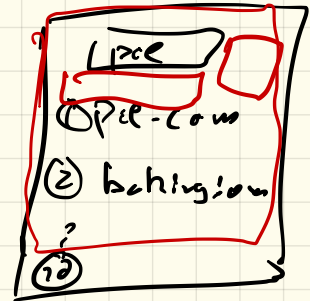
inverted index

key word



Big challenges

• Ranking pages for key words



Rank web page on keyword?

→ pie

k-grams on text on page
+ Jaccard w/ {pie}
+ Minhashing

All
Visits

cosine similarity

$v_{pie} = (0, 0, 0, \dots, 0, 1)$

bow page (1, 2, ~~4~~, ..., ^{pie}7, 1, 0)

What can go wrong?

↳ many copies of "pie"

How did search engines know about pages?

→ Crawlers: goes to webpage, follows links, & records page.

hyperlinks $\langle a \ href = " \text{newpage} " \rangle \text{ text } \langle / a \rangle$

↑
very informative

words on page

↳ put into how rep for page

page how vec (v_1, v_2, \dots, v_n)

text hyperlinks $(v_1', v_2', \dots, v_n')$
 $(\bar{v}_1, \bar{v}_2, \dots, \bar{v}_n)$

Index

pages hand-curated list of links

↳ Yahoo!

Look Smart

1. Google
2. Youtube
3. Facebook
4. Baidu
5. Wikipedia
6. Reddit
7. Yahoo.
8. Google India
9. Tencent QQ
10. Amazon

Page Rank

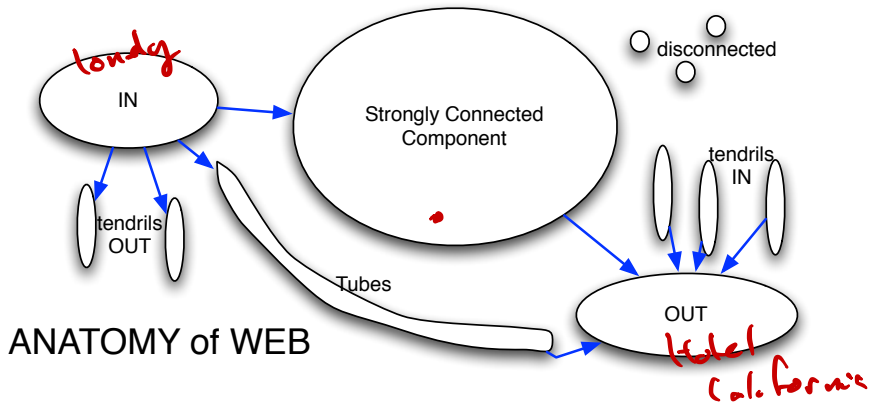
Idea #1, Pages are important if
linked to by other important pages.

Idea #2 How likely a random surfer
would find this page.

Markov Chain \leftarrow model of web graph.

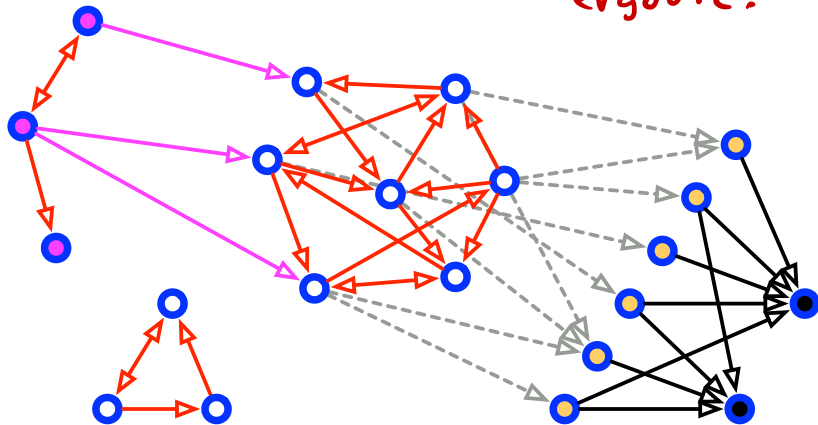
ergodic?

Anatomy of Web



Anatomy of Web

Not ergodic!



Teleportation (Taxation)

Idea 15% steps, jump to random page.

$$g_x = P^n g_0$$

$P_{n \times n}$

$$\beta = 0.15$$

$$g_{i+1} = \left((1-\beta)P + \beta Q \right) g_i$$

Q all $1/m$

P
ergodic

$$g_x^i = (P^i)^n g_0$$

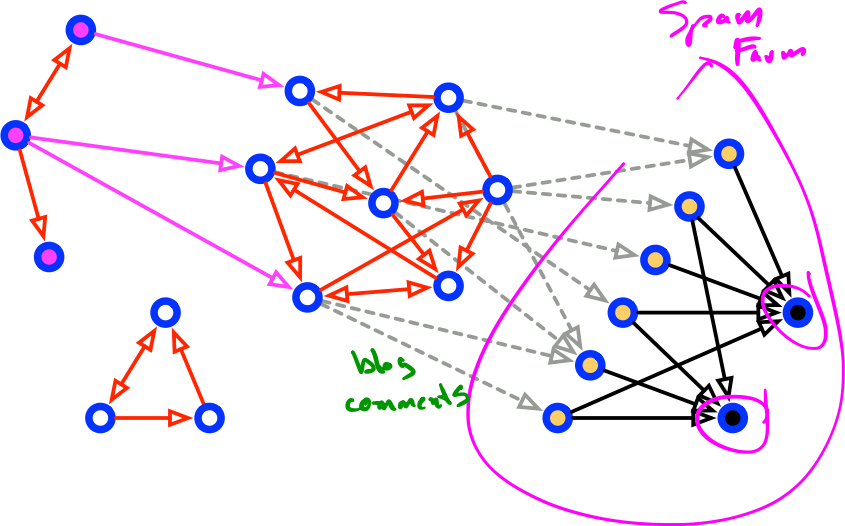
Page Ranks vec

$$\begin{pmatrix} 0.15 \\ 0.15 \\ \vdots \\ 0.15 \\ m \end{pmatrix}$$

Data Matrix

0
0
1/100
1/100
1/100
0
1/100
0
...

Spam Farms



$$\text{eigs}(M) = [V, L]$$

$$L = \begin{pmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \dots \\ & & & \lambda_n \end{pmatrix}$$

$$M = V L V^T$$

$$\sqrt{M} = V \sqrt{L}$$

$$\sqrt{L} = \begin{pmatrix} \sqrt{\lambda_1} & & \\ & \sqrt{\lambda_2} & \\ & & \dots \\ & & & \sqrt{\lambda_n} \end{pmatrix}$$