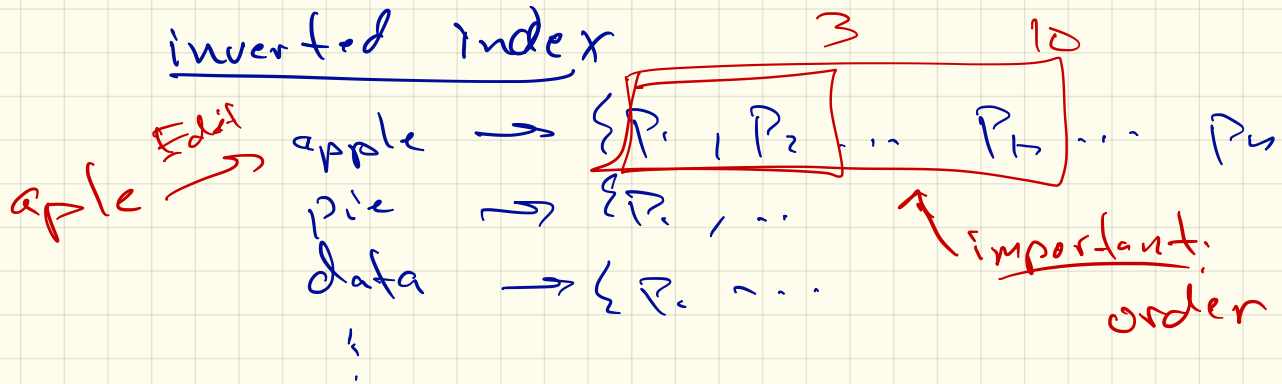


# Page Rank

## and Search Engines

How does a search engine work?



# Web page Similarity

pre 1998  
PR-PageRank

Altavista, Lycos, Infoseek

web page  $\rightarrow$  html text  $\rightarrow$  bag-of-words

$$P_i \rightarrow v_i \in \mathbb{R}^{10,000} \quad v_i = (0, 0, 8, 0, 0, \dots, 0, 5, 0)$$

apple pie

search "apple pie"  $\downarrow$

$$q \in \mathbb{R}^{10,000} \quad q_i = (0, \dots, \frac{1}{\sqrt{2}}, 0, \dots, \frac{1}{\sqrt{2}}, 0)$$

highest cosine-sim( $v_i, q$ )  $\rightarrow$  top of list.

... but? "apple pie", apple pie, ... white

Wattle : search engine vs. spammer's

---

modify dist  
 $ave(cos, J_{accasd})$   
emphasize certain words  
cap word count.

query: ( o o o ~~copy~~ )  
modify

copy top  
pages into  
bottom of  
your  
page

# Index

Yahoo! / Look Smart

business model : Paid placement

1. Google
2. Youtube
3. Facebook
4. Baidu (China)
5. Wizi
6. (18) Tencent QQ
7. Teobao
8. Tmall
9. (6) : Yahoo!
10. (11) Amazon
11. (7) Twitter

15 (10) Instagram

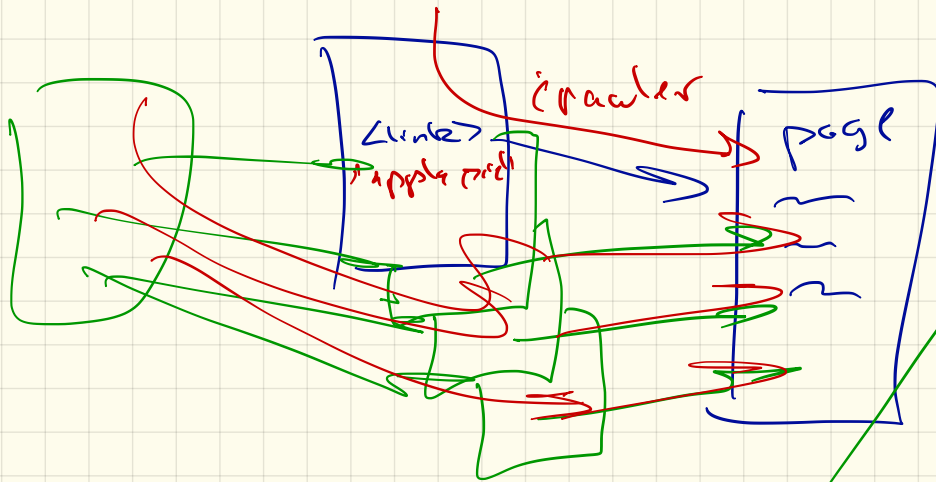
23 (17) Netflix

# Crawlers

automated bot

↳ moves around the web

→ follow links



...but

# Page Ranks

Idea 1 important webpages link to important webpages

Idea 2 Important webpages are visited often by a random surfer

Markov Chain



"crawler"

Model web as graph  $G = (V, E)$

Vertices = pages  
Edges = hyperlinks

web graph

$G = (V, E) \rightarrow$  Adjacency matrix

Probability  $\leftarrow A$

Transition Matrix  $P$

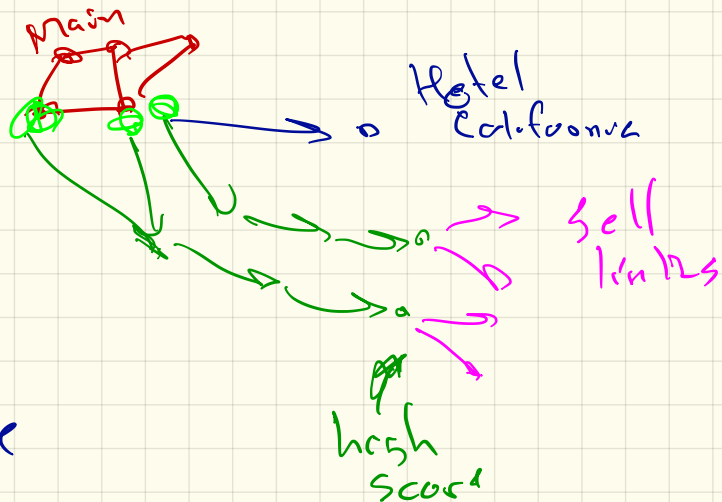
$r: v_i = \text{eig}(P, 1) \equiv$  Page rank vector

$r(\text{page})$ : larger better

Score (page, query) = magic  $\left( \overset{\text{page rank}}{\boxed{r(\text{page})}}, \underset{\substack{\cos(\theta_{p,q}) \\ \text{hyperlinks}}}{r} \right)$

# Is the webgraph ergodic?

- cycles? No
- connected? No (bot ok)
- transient?



Fix: Teleportation

$\beta = 15\%$  jump to random page



Compute

$$g^* = P^* g_0$$

option 1

eigs (P) <sup>house</sup>

option 2

compute  $P^n$  large  $n$

for any  $g_0$ :  $g_x \approx P^n g_0$

"small world"  $P^7$  <sup>vers</sup> dense

option 3

for  $i = 1 \dots (n=50)$

$$g_{i+1} = P g_i$$

return  $g_x = g_n$

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

$$Q = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \frac{1}{n}$$

$$= (1-\beta)P g_i + \beta \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$$

random jump

2015: "truth"

# Trust Rank

Trust more: Wikipedia, .edu domain

Run 2 versions of PageRank.

• Regular: teleportation is uniform  
 $r(p)$

• Trust: teleportation is more likely  
to jump to trusted page.  
 $t(p)$

$$S(p) = \frac{r(p) - t(p)}{r(p)} : \text{larger} \rightarrow \text{more likely spam}$$