relevant 相关的
Converge 收敛

Sparse matrix 稀疏矩阵

1. dead end

no out

2. spider trap

只在一个group里, 不能出group

3. teleport : trap 和 dead end 执行策略.

∠

# Some problems with Page Rank:

① 受特定主题影响：Topic-Specific Page Rank

② 单一衡量：Hubs-and-Authorities 枢纽和权威

③ 垃圾链接：Trust Rank
Spam

Topic-Specific Page Rank    Aallow be answered based
准备一个相关 set. teleport 对象在此里找    user's interests

计算方式：只需改变A.   $A_{ij} = \begin{cases} \beta M_{ij} + (1-\beta)/|S| & \text{if } i \in S \\ \beta M_{ij} + 0 & \text{else} \end{cases}$

关于选择 topic：
$\begin{cases} 可以让 user 从菜单选 \\ 可以 将查询往主题夹 \\ 可以 使用查询上夹. \end{cases}$

# Proximity in Graphs 图中接近度
短路径并不一定好.

# Spam: web pages that are the result of any deliberate action to boost a web page's position in search engine results, incommensurate with page's real value.

早期 Page search: crawel the web, index pages 包含查询关键字.
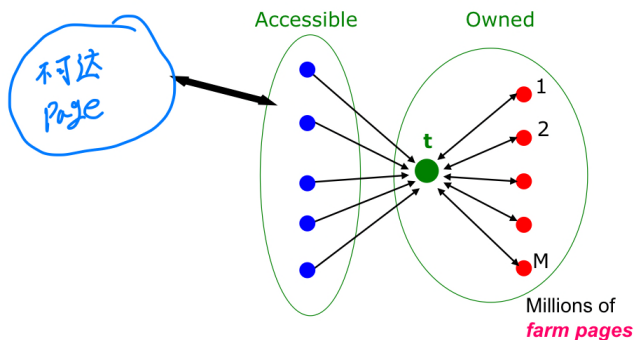
早期 Page Rank: ⎰ 关键字出现的次数.
⎱ 关键字出现的位置, eg: title

<span style="color:blue">锚文本</span>

对于 early spam: ⎰ 用 anchor text 或 link 周围文本而非 link 自说的文本.
⎱ PageRank 可以帮助过滤掉 spam.

Round 2 - spam farm: concentrate PageRank on a single page.

↳实现: ① link spam: create link structures that boost pageRank of a particular page

从 accessible pages <span style="color:blue">(如博客评论区)</span> 尽可能多创建指向 target page t 的 Link.



**Accessible** **Own**

不可达 page

t

<span style="color:pink">经典组织图.

N pages on the web

M pages spammer owned

X: PageRank by accessible

Y: PageRank of t.</span>

Rank of each form page: $\frac{\beta y}{M} + \frac{1-\beta}{N}$

而 $y = x + \beta M [\frac{\beta y}{M} + \frac{1-\beta}{N}] + \frac{1-\beta}{N}$

$= x + \beta^2 y + \frac{\beta(1-\beta)M}{N} + \frac{1-\beta}{N}$ <span style="color:pink">too small so ignore ✕</span>

SO: $y = \frac{x}{1-\beta^2} + c\frac{M}{N}$ 其中 $c = \frac{\beta}{1+\beta}$ <span style="color:pink">(can make M large to make PageRank large)</span>

Combating term spam:

- analyze text using statistical methods similar to email spam filtering.
- Detecting approximate duplicate pages.

Combating link spam:

Detection and blacklisting of structures that look like slam farms.

TRUST Rank: topic-specific PageRank with a teleport set of trusted **P30**

原理: 近似隔离 approximate isolation
good pages rardy to point to be pages

每顶面信任度在 0-1, 信任度是加法, 来自它的 inlinks 传送的信任度之和.

eg: P的 trust 是 tp, P has a set of out-links Op. seed 设为1.
P 传送给它指向的 q 的 trust 为 $\frac{\beta t_p}{|O_p|}$

if teleport set 都是 trusted pages, 则 TrustRank = PageRank

信任削减

- Trust attenuation: 路径越长, 越削减.
- Trust splitting: split 越多, 传送的 trust 越少.

→ 信任传播模型. → solution 1

How to Pick a seed set of K pages:

① PageRank: 选 K↑ PageRank最高的. 因为现实中很难让 bad page 有高PageRank

② Use trusted domains. eg: .edu .mil .gov.


Solution2:  Spam Mass Estimate

$r_p$ = Page P 的 PageRank

$r_p^+$ = PageRank of P的从 trusted Pages teleport Into 的.

from spam pages $r_p^- = r_p - r_p^+$

Spam mass of P = $\dfrac{r_p^-}{r_p^+}$


Hubs and Authorities:

HITS: Hypertext - Induced Topic Selection

评估 Pages or documents 重要程度的方法.

→ 每1 page 有2个 scores

中心, expert

① as hub: 指向的 authority 的 sum 和 (value)
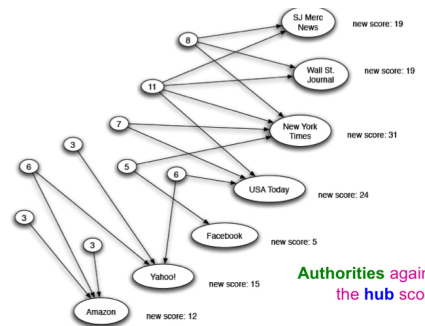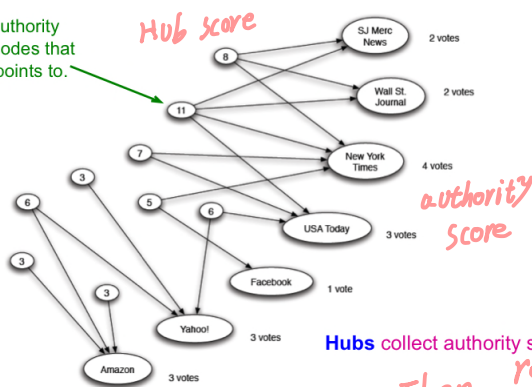
② as authority: 来自 hub 的 vote总和

content

Authority: Pages 包含有用信息              hub: 指向authority 的 Pages

Like 报纸homepage                         like 报纸列表

Sum of authority scores of nodes that the node points to.

*Hub score*

*authority score*

1 vote

3 votes

3 votes

*Hubs colle*

*Then* *reweight* →

new score: 5

new score: 15

new score: 12

*Autho...*

# HITS 算法流程:

初始化: $a_j^{(0)} = \dfrac{1}{\sqrt{N}}$    $h_j^{(0)} = \dfrac{1}{\sqrt{N}}$

$a_t^{(t+1)} = \sum\limits_{j \to i} h_j^{(t)}$   then   $h_i^{(t+1)} = \sum\limits_{i \to j} a_j^{(t)}$

*就是变成*

$a_i^{(t+1)} = \dfrac{a_i^{(t)}}{\sqrt{\sum\limits_{i}(a_i^{(t)})^2}}$

the 归一化: $\sum\limits_i (a_i^{(t+1)})^2 = 1$,   $\sum\limits_j (h_j^{(t+1)})^2 = 1$

until 收敛.

## HITS 矩阵计算:

邻接矩阵 $A_{N \times N}$, 其中 $A_{ij} = 1$    if $i \to j$
$A_{ij} = 0$    else

故 $h_i = \sum\limits_{i \to j} a_j$ 可以写成 $h_i = \sum\limits_j A_{ij} \cdot a_j$ 即 $h = A \cdot a$

同理  $a_i = \sum\limits_{j \to i} h_j$ 可以写成 $a_i = \sum\limits_j A_{ji} \cdot h_j$ 即 $a = A^T \cdot h$

故

st1:  set $a_i = h_i = \dfrac{1}{\sqrt{n}}$

st2: Repeat $\begin{cases} h = A \cdot a \\ a = A^T \cdot h \\ a \, 与 \, h \, 归一化 \end{cases}$  Until $\begin{cases} \sum\limits_i (h_i^{(t)} - h_i^{(t-1)})^2 < \varepsilon \\ \sum\limits_i (a_i^{(t)} - a_i^{(t-1)})^2 < \varepsilon \end{cases}$

Then $a = A^T \cdot (A \cdot a)$    故    $\begin{cases} a = A^T(Aa) = (A^TA)a \\ h = A(A^Th) = (AA^T)h \end{cases}$
$\underset{new\ h}{\underbrace{}}$    更新
      成

from $U$ to $V$ 中

Summary: PageRank 与 HITS 想解决的问题都隐 $U$-$V$ 的 in-link 价值.

$\begin{cases} \text{PageRank: 取决于 links to } u. \\ \text{HITS: \quad 取决于 links out of } u. \end{cases}$