

Web 分析

- コンテンツ
 - － 要約, トピック検出
- リンク構造
 - － ページの人気度, Webコミュニティの発見
- 利用履歴, アクセスログ
 - － そのサイトのアクセスログ, 検索エンジンでのアクセスログ
 - － ページの人気度, サイトデザインの改善
- 更新履歴
 - － クローラのスケジューリング

1

Web 検索エンジンにおけるランキング手法

- 初期の検索エンジンのランキング手法
 - － 「検索式に強くマッチするページ \equiv ユーザが求めるページ」
 - － 例: tf-idf 法, 特定のタグ内の単語
- 90年代半ばより Web では
 - － 「検索式に強くマッチするページ \neq 有用なページ」
 - － 個人の日記ページ
 - － SPAM: ページをランキング上位に押し上げるための小細工
例: よく検索されるキーワードを小さなフォントや見えない色のフォントでたくさん埋め込む
- 解決策: 「他ページからリンクされるページ \equiv 有用なページ」
⇒ **リンクの情報を利用したランキング手法**

2

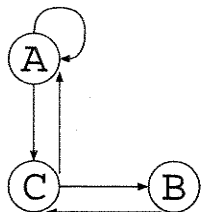
PageRankTM (google.com)**概要**

- 多くのページからリンクされるページは有用なページだろう
 - 多くの有用なページからリンクされていればなお有用なページ
- ↓
- 各ページに「有用度」を割り当てる
 - 有用度をリンクに沿ってページからページへと波及させる
 - ページから複数のリンクが出ている場合, そのページの有用度をそれらへ均等に分割

3

PageRank™ (google.com)

計算方法



$$\begin{pmatrix} A_0 \\ B_0 \\ C_0 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

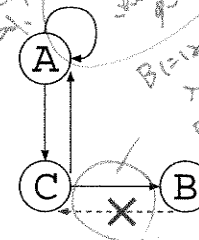
$$\begin{pmatrix} A_{i+1} \\ B_{i+1} \\ C_{i+1} \end{pmatrix} = \begin{pmatrix} 1/2 & 0 & 1/2 \\ 0 & 0 & 1/2 \\ 1/2 & 1 & 0 \end{pmatrix} \begin{pmatrix} A_i \\ B_i \\ C_i \end{pmatrix}$$

$$\lim_{i \rightarrow \infty} \begin{pmatrix} A_i \\ B_i \\ C_i \end{pmatrix} = \begin{pmatrix} 6/5 \\ 3/5 \\ 6/5 \end{pmatrix} \leftarrow \dots \leftarrow \begin{pmatrix} 9/8 \\ 1/2 \\ 11/8 \end{pmatrix} \leftarrow \begin{pmatrix} 5/4 \\ 3/4 \\ 1 \end{pmatrix} \leftarrow \begin{pmatrix} 1 \\ 1/2 \\ 3/2 \end{pmatrix} \leftarrow \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

4

PageRank™ (google.com)

dead end problem



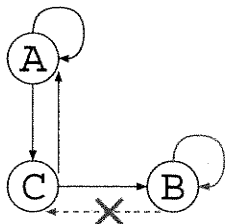
$$\begin{pmatrix} A_{i+1} \\ B_{i+1} \\ C_{i+1} \end{pmatrix} = \begin{pmatrix} 1/2 & 0 & 1/2 \\ 0 & 0 & 1/2 \\ 1/2 & 0 & 0 \end{pmatrix} \begin{pmatrix} A_i \\ B_i \\ C_i \end{pmatrix}$$

$$\lim_{i \rightarrow \infty} \begin{pmatrix} A_i \\ B_i \\ C_i \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \leftarrow \dots \leftarrow \begin{pmatrix} 5/8 \\ 1/4 \\ 3/8 \end{pmatrix} \leftarrow \begin{pmatrix} 3/4 \\ 1/4 \\ 1/2 \end{pmatrix} \leftarrow \begin{pmatrix} 1 \\ 1/2 \\ 1/2 \end{pmatrix} \leftarrow \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

5

PageRank™ (google.com)

spider trap problem



$$\begin{pmatrix} A_{i+1} \\ B_{i+1} \\ C_{i+1} \end{pmatrix} = \begin{pmatrix} 1/2 & 0 & 1/2 \\ 0 & 1 & 1/2 \\ 1/2 & 0 & 0 \end{pmatrix} \begin{pmatrix} A_i \\ B_i \\ C_i \end{pmatrix}$$

$$\lim_{i \rightarrow \infty} \begin{pmatrix} A_i \\ B_i \\ C_i \end{pmatrix} = \begin{pmatrix} 0 \\ 3 \\ 0 \end{pmatrix} \leftarrow \dots \leftarrow \begin{pmatrix} 5/8 \\ 2 \\ 3/8 \end{pmatrix} \leftarrow \begin{pmatrix} 3/4 \\ 7/4 \\ 1/2 \end{pmatrix} \leftarrow \begin{pmatrix} 1 \\ 3/2 \\ 1/2 \end{pmatrix} \leftarrow \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

6

PageRank™ (google.com)

dead end problem と spider trap problem の回避

$$\begin{pmatrix} A_{i+1} \\ B_{i+1} \\ C_{i+1} \end{pmatrix} = (1 - tax) \begin{pmatrix} 1/2 & 0 & 1/2 \\ 0 & 1 & 1/2 \\ 1/2 & 0 & 0 \end{pmatrix} \begin{pmatrix} A_i \\ B_i \\ C_i \end{pmatrix} + tax \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

↑
spider trap を
避けるために税を徴収

↑
dead end problem を
避けるために税を再分配

$$\lim_{i \rightarrow \infty} \begin{pmatrix} A_i \\ B_i \\ C_i \end{pmatrix} = \begin{pmatrix} 7/11 \\ 21/11 \\ 5/11 \end{pmatrix} \leftarrow \dots \leftarrow \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \quad (tax = 0.2)$$

7

PageRank™ (google.com)

ランダムウォークとしての解釈

$$\begin{pmatrix} A_{i+1} \\ B_{i+1} \\ C_{i+1} \end{pmatrix} = (1 - \epsilon) \begin{pmatrix} 1/2 & 0 & 1/2 \\ 0 & 1 & 1/2 \\ 1/2 & 0 & 0 \end{pmatrix} \begin{pmatrix} A_i \\ B_i \\ C_i \end{pmatrix} + \epsilon \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

\uparrow \uparrow \uparrow \uparrow
 時間 $i+1$ に ジャンプ 時間 i に ランダムに
 各ページに しない 各ページに ジャンプ
 いる確率 場合 いる確率 する場合

ϵ = 「リンクをたどらずにランダムにジャンプする確率」

8

ハブーオーソリティ解析

概要

- ページには二種類の情報がある
 - 情報を提供するページ
 - 情報を持つページへのリンクを提供するページ
(例: ポータル, リンク集ページ)
- オーソリティ = 有用な情報を持つページ
- ハブ = 有用なリンクを持つページ

HITS (Hypertext Induced Topic Search)

“Authoritative Sources in a Hyperlinked Environment”

by Jon M. Kleinberg, Journal of ACM, 46(5), 1999

- 与えられたキーワードを含むページ集合を取得
(Focused Subgraph)
- このページ集合に対してハブーオーソリティ解析
そのキーワードに対応するトピックに関する「オーソリティ」と「ハブ」を求める.
- YST (Yahoo! Search Technology) が採用したと言われている.
- Jon Kleinberg が 2006 年, IMU Nevanlinna 賞を受賞

9

ハブーオーソリティ解析

概要

各ページに「オーソリティ度」と「ハブ度」を割り当てる

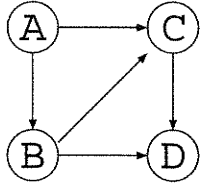
- ページのオーソリティ度とは?
多くの良いハブからリンクされているページは良いオーソリティ
- ページのハブ度とは?
多くの良いオーソリティをリンクしているページは良いハブ

相互再帰的に定義される

ハブーオーソリティ解析

計算方法

$$\vec{h}_{i+1} = A \vec{a}_i$$



$$\begin{pmatrix} 2 \\ 2 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

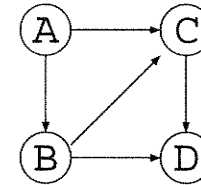
\uparrow \uparrow \uparrow
 ハブ度 隣接行列 A オーソリティ度

12

ハブーオーソリティ解析

計算方法

$$\vec{a}_{i+1} = A^t \vec{h}_i$$



$$\begin{pmatrix} 0 \\ 2 \\ 4 \\ 3 \end{pmatrix} = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 \end{pmatrix} \begin{pmatrix} 2 \\ 2 \\ 1 \\ 0 \end{pmatrix}$$

\uparrow \uparrow \uparrow
 オーソリティ度 A^t ハブ度

13

ハブーオーソリティ解析

計算方法

$$A^t A = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 \end{pmatrix} \begin{pmatrix} 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{pmatrix} = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 1 & 2 & 1 \\ 0 & 0 & 1 & 2 \end{pmatrix}$$

オーソリティ度: $\vec{a}_{2i} = (A^t A)^i \vec{a}_0$

$$(A^t A)^i \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} \rightarrow \begin{pmatrix} 0 \\ 2 \\ 4 \\ 3 \end{pmatrix} \rightarrow \begin{pmatrix} 0 \\ 6 \\ 13 \\ 10 \end{pmatrix} \rightarrow \begin{pmatrix} 0 \\ 19 \\ 42 \\ 33 \end{pmatrix} \rightarrow \begin{pmatrix} 0 \\ 61 \\ 136 \\ 108 \end{pmatrix} \rightarrow \dots$$

14

ハブーオーソリティ解析

計算方法

$A^t A$ は実数対称行列

→ 以下を満たす対角行列 D が存在

$$A^t A = U^{-1} D U$$

よって

$$(A^t A)^i \vec{a}_0 = U^{-1} D^i U \vec{a}_0$$

通常は、さらに正規化を行う

$$\vec{a}_{i+2} = \frac{A^t A \vec{a}_i}{|A^t A \vec{a}_i|}$$

15

ハブーオーソリティ解析

計算方法

ハブ度も同様:

$$\vec{h}_{2i} = (AA^t)^i \vec{h}_0$$

AA^t は実数対称行列

→ 以下を満たす対角行列 D_h が存在

$$AA^t = U_h^{-1} D_h U_h$$

よって

$$(AA^t)^i \vec{h}_0 = U_h^{-1} D_h^i U_h \vec{h}_0$$

正規化:

$$\vec{a}_{i+2} = \frac{AA^t \vec{a}_i}{|AA^t \vec{a}_i|}$$

16

ハブーオーソリティ解析

リンク関係の主成分分析としての解釈

- $AA^t(i, j)$ = ノード i, j が共にリンクしているノード数
- $A^t A(i, j)$ = ノード i, j を共にリンクしているノード数

↓

- h は「どんなノードをリンクしているか」に関する主成分分析の第一主成分の係数
- a は「どんなノードからリンクされているか」に関する主成分分析の第一主成分の係数

18

ハブーオーソリティ解析

リンク関係の主成分分析としての解釈

$$\begin{cases} h = \lim_{i \rightarrow \infty} h_i \\ a = \lim_{i \rightarrow \infty} a_i \end{cases} \text{ とおくと } \begin{cases} h = \mu Aa \\ a = \eta A^t h \end{cases} \text{ よって } \begin{cases} h = \mu \eta AA^t h \\ a = \mu \eta A^t Aa \end{cases}$$

$\mu \eta = \lambda^{-1}$ とおくと,

$$\begin{cases} \lambda h = AA^t h \\ \lambda a = A^t Aa \end{cases}$$

h, a は $AA^t, A^t A$ の固有ベクトル
(最大の固有値に対する)

17

