

Patterns and models that are?

valid useful unexpected understandable

Data Mining tasks.

- ①. Descriptive methods eg. cluster
- ②. Predictive methods eg. recommender sys

Usage Quality context screaming
scalability 扩展

Links as votes. ① more links, Page more impor

② Links from important Page count more

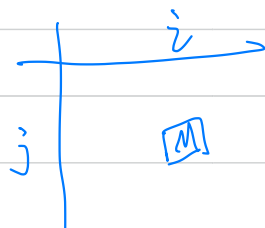
in-link: 点链接至该页面 $i \rightarrow j$

Value j = 来源网页分支

$$\sum_i v_{ij} = 1$$

$r = M \cdot r$

All Page score



$M_{ji} = \frac{1}{n_i}$

eigenvector: 特征向量

$$(Ax = \lambda x) \text{ why} \quad \text{其中 } A \text{ 是 } N \text{ 矩阵}$$

$$y^{(t+1)} = \begin{matrix} \boxed{} & \boxed{} \\ n \times n & n \times 1 \\ M & \underline{y^{(t)}} \end{matrix}$$

Link analysis approaches: { Page Rank
Topic-specific Page Rank
Web spam detecting algorithms

Page Rank: { 有更多link的 page 更重要. in-link like vote

Link from 重要页面的 count more.

公式描述: { - 1 page 的 score 为它所有 in-link 权重之和 $\rightarrow r_j$
[每个 link 的权重等于它 source page 权重 / out-link 数目.

$$r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i} \quad \text{其中 } d_i \text{ 为节点 } i \text{ 的出度数目}$$



$$M_{ji} = \frac{1}{d_i} \quad M_{ii} = 0$$

\swarrow M 列的 sum = 1.

M_i 列随机矩阵

权重向量 r : 代表每个 page 的 score

$$\sum_i r_i = 1 \quad \leftarrow \text{所有页面 score 之和为 1}$$

$$M \cdot r = r$$

特征向量: $AX = \lambda X$ 则 X 是特征向量, λ 是特征值. 故 1 是特征值
考虑到 r 向量中每个 r_i 之和为 1 而 M 的每个 column 之和为 1 故 $Mr \leq 1$

M_{ji} 指从 i 出发到 j 的值是 $\frac{1}{d_i}$

	出度	0	1	2	j
目的 i	0				
	1				
	2				

迭代法: 依据 $r = Mr$ 计算

初始化: $r^{(0)} = \begin{bmatrix} \frac{1}{N} \\ \vdots \\ \frac{1}{N} \end{bmatrix}$

N 为结点数目

$$r^{(t+1)} = M \cdot r^{(t)} \quad \text{until } |r^{(t+1)} - r^{(t)}| < \epsilon$$

$$\text{总结 } r^{(n)} = M \cdot r^{(n-1)} = M^n r^{(0)}$$

结果接近主特征向量

β 随机游走: $P(t+1) = M \cdot P(t)$, 假设达到一个状态 $P(t+1) = M \cdot P(t) = P(t)$
则 $P(t)$ 是 stationary distribution

结论: 对于满足特定条件的图形, 静态分布是唯一的, 无论 $t=0$ 时初始概率为何, 最终都达到静态分布.

Each time step, 都执行

Page Rank 的问题是: 是 problem

① dead ends: 结点无出点 to go out

r 很易变成形如 $\begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$

② spider traps: 局限在某 group.

trap 很易导致最终的 r 为类似 $\begin{bmatrix} 0 \\ \vdots \\ 1 \end{bmatrix}$

不算 problem.

Trap 的 solution: Teleports.

每次跳有 2 个 choice $\begin{cases} \beta \text{ 概率 follow link random} \\ 1-\beta \text{ Prob: jump to random page} \end{cases}$

β 一般 $0.8 \sim 0.9$

Dead End 的 solution: TelePort too!

1.0 概率随机跳

sparse matrix: 稀疏矩阵

使用 TelePort 之后, 公式变成 $r_j = \sum_{i \rightarrow j} \beta \frac{r_i}{d_i} + (1-\beta) \frac{1}{N}$

即 $A = \beta M + (1-\beta) \left[\frac{1}{N} \right]_{N \times N}$

$r = Ar = \beta Mr + \underbrace{\left[\frac{1-\beta}{N} \right]_{N \times N}}_{\text{为常数, 故}} \rightarrow \text{since } \sum_i r_i = 1$

直后加 $\frac{1-\beta}{N}$

