

Dijkstra 算法在 Web 结构挖掘的应用

林耀进¹, 王晨曦²

(1、漳州师范学院计算机科学与工程系 福建 漳州 353000 2、漳州职业技术学院计算机工程系 福建 漳州 353000)

摘要】 该文从 Web 结构挖掘角度出发, 利用概率论分析了 Web 结构挖掘的 PageRank 算法, 得出挖掘结果, 最后介绍 Dijkstra 算法在其挖掘结果的应用。

关键词】 Web 结构挖掘, PageRank 算法, Dijkstra 算法, 权重

搜索引擎 Google 的成功, 取决于它采用了有效的 Web 信息挖掘技术。Web 挖掘指在 WWW 上挖掘潜在的、有用的模式及隐藏的信息过程^[1]。Web 挖掘分为 Web 内容挖掘、Web 结构挖掘和 Web 使用挖掘。其中结构挖掘则是从人为地链接结构中获取有用知识的过程。在设计搜索引擎等服务时, 对 Web 页面的连接结构进行挖掘以得出有用的知识是提高检索的重要手段。

本文从 Web 结构挖掘入手, 对 Web 结构挖掘的 PageRank 算法结合概率论进行研究应用。然后介绍 Dijkstra 算法在挖掘结果中的应用。

1. 算法的介绍

1.1 PageRank 算法

在 PageRank 方法中的 PageRank 被定义为^[2]: 设 u 是一个 Web 页, F_u 为所有 u 指向的页面的集合, B_u 为所有指向 u 的页面的集合。设 $N_u = |F_u|$ 为从 u 发出的链接的个数, 那么 u 页面的 PageRank 可以定义为:

$$R(u) = c \sum_{v \in B_u} R(v) / N_v \quad (1)$$

其中 $c(1)$ 为归一化因子(因为所有页面的 RankPage 之和为一个常数)。PageRank 算法的实现过程: 将网页的 URL 对应成唯一的整数, 把每一个超链接用其整数 ID 存放到索引数据库中, 经过预处理后, 设每个网页的初始 PR 值为 1, 通过以上的递归算法计算每一个网页的 PageRank 值, 反复进行迭代, 直至结构收敛。显然, PageRank 值越大, 该页面权威性越高。

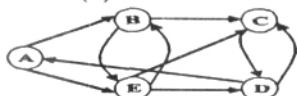
1.2 Dijkstra 算法

Dijkstra 算法是由荷兰计算机科学家艾兹格·迪科斯彻发现的。Dijkstra 算法是图论学中求解最短路径问题的经典算法, Dijkstra 算法建立在抽象的网络模型上, 把路抽象为网络中的边, 以边的权值来表示路相关的参数, 算法确定了赋权网络中从某点到所有其它结点的具有最小权的路。权的含义是广泛的, 可以表示距离、费用、数量等等^[3]。

Dijkstra 算法的输入包含了一个有权重的有向图 G , 以及 G 中的一个来源顶点 S 。我们以 V 表示 G 中所有顶点的集合。每一个图中的边, 都是两个顶点所形成的有序元素对。 (u, v) 表示从顶点 u 到 v 有路径相连。我们以 E 所有边的集合, 而边的权重则由权重函数 $w: E \rightarrow [0, \infty]$ 定义。因此, $w(u, v)$ 就是从顶点 u 到顶点 v 的非负花费值 (cost)。边的花费可以想像成两个顶点之间的距离。任两点间路径的花费值, 就是该路径上所有边的花费值总和。已知有 V 中有顶点 s 及 t , Dijkstra 算法可以找到 s 到 t 的最低花费路径 (i.e. 最短路径)。这个算法也可以在一个图中, 找到从一个顶点 s 到任何其他顶点的最短路径。

2. 仿真实验

2.1 根据图(1)的 Web 页面结构和公式(1)的 PageRank 定义, 很容易列出 Web 的 PageRank 之间的线性关系, 不妨引入记号 $u=R(u)$ 则有公式(2)



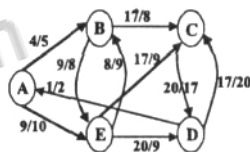
图(1) Web 页面结构

对公式(2)求解得公式(3)

$$\begin{cases} A = \frac{D}{2} \\ B = \frac{A}{2} + \frac{E}{3} \\ C = \frac{B}{2} + \frac{D}{2} + \frac{E}{3} \\ D = C + \frac{E}{3} \\ E = \frac{A}{2} + \frac{B}{2} \end{cases} \quad (2)$$

$$\begin{cases} 2A = D \\ 9A = 10E \\ 4A = 5B \\ 17B = 8C \\ 9B = 8E \\ 20C = 17D \\ 9C = 17E \\ 20E = 9D \end{cases} \quad (3)$$

对于页面 B, E, 其中页面 B, E 存在直接链接关系: 把 $W(B, E)$ 声明为代表 $B \rightarrow E$ 的回顾因子权重。 $W(,)$ 表示两个页面之间的紧密链接关系, $W(,)$ 值越大说明两个页面之间关系越紧密。如 $W(B, E) = 9/8$, $W(E, B) = 8/9$ 。根据图(1) Web 页面结构和表达式(3), 很容易得出图(2) Web 带权重的页面结构。



图(2) Web 带权重的页面结构

2.2 在图(2)中, 我们发现对于每一对页面 P_i, P_j (P 代表页面的集合), P_i, P_j 若从 P_i 到 P_j 存在路径数 2 , $W(P_i, P_j)$ 都是一样的。如对图(2), 页面 A 到 D 存在多条路径, 其中: $R_1: A \rightarrow B \rightarrow C \rightarrow D$; $R_2: A \rightarrow E \rightarrow D$; $R_3: A \rightarrow B \rightarrow E \rightarrow D$; $R_4: A \rightarrow E \rightarrow B \rightarrow C \rightarrow D$ 。 $W(A, D)$ 都为 2。依据概率论的积事件, 说明若从页面 A 搜索页面 D, 无论走哪条路径, 概率是等同, 即 $P(A \rightarrow B \rightarrow C \rightarrow D) = P(A \rightarrow E \rightarrow D) = P(A \rightarrow B \rightarrow E \rightarrow D) = P(A \rightarrow E \rightarrow B \rightarrow C \rightarrow D)$ 。

同样, 在图(2)中, 我们发现各条路径的回顾因子权重和不一定相等。如路径 $R_1: W_{R1} = 4/5 + 17/8 + 20/17 = 2789/680$; 路径 $R_2: W_{R2} = 9/10 + 20/9 = 281/90$; 路径 $R_3: W_{R3} = 4/5 + 9/8 + 20/9 = 1259/360$; 路径 $R_4: W_{R4} = 9/10 + 8/9 + 17/8 + 20/17 = 62306/12240$ 。即 $W_{R1} < W_{R2} < W_{R3} < W_{R4}$ 其中 W_{Ri} 代表路径 R_i 的权重和。依据概率论的和事件, W_{Ri} 越大说明, 路径 R_i 的页面出现的概率越大, 如在图(2)中, 当从页面 A 搜索页面 D 时, 路径 R_4 中的页面出现概率较大。

2.3 根据前面推出的结果, 利用 Dijkstra 算法对带回顾因子权重的 Web 页面结构图求两点之间的最短路径, 其伪代码如下 (其中 u, v 为顶点, $d[v]$ 中储存的便是从 s 到 v 的最短路径, $w(u, v)$ 代表顶点 u 与顶点 v 之间的权重):

```

1 function Dijkstra(G, w, s)
2   for each vertex v in V[G] // 初始化
3     d[v] := infinity
4     previous[v] := undefined
5   d[s] := 0
6   S := empty set
7   Q := set of all vertices
8   while Q is not an empty set // Dijkstra 算法主体
9     u := Extract_Min(Q)
10    S := S union {u}
11    for each edge (u, v) outgoing from u
12      if d[v] > d[u] + w(u, v) // 拓展边(u, v)
13        d[v] := d[u] + w(u, v)

```

(下转第 103 页)



论文写作，论文降重，
论文格式排版，论文发表，
专业硕博团队，十年论文服务经验



SCI期刊发表，论文润色，
英文翻译，提供全流程发表支持
全程美籍资深编辑顾问贴心服务

免费论文查重：<http://free.paperyy.com>

3亿免费文献下载：<http://www.ixueshu.com>

超值论文自动降重：http://www.paperyy.com/reduce_repetition

PPT免费模版下载：<http://ppt.ixueshu.com>

阅读此文的还阅读了：

- [1. HITS算法在Web挖掘中的应用与改进](#)
- [2. 一种基于云计算的Web结构挖掘算法](#)
- [3. Web结构挖掘算法探讨](#)
- [4. 基于Dijkstra算法的Web服务合成选择策略研究](#)
- [5. 专业网站零输入导航引擎设计与实现](#)
- [6. Dijkstra算法在智能公交查询系统中的应用](#)
- [7. Dijkstra改进算法在地震救援中的应用](#)
- [8. Dijkstra改进算法在机器人避障问题的应用](#)
- [9. Web数据挖掘算法](#)
- [10. Dijkstra算法在Web结构挖掘的应用](#)
- [11. 基于 Dijkstra算法的](#)
- [12. 结构挖掘中web有向图模型的改进算法](#)
- [13. Trawling算法在Web结构挖掘中的应用](#)
- [14. Web结构挖掘及其算法分析](#)
- [15. Web结构挖掘算法研究](#)
- [16. Dijkstra算法在最优投资策略问题中的应用](#)

[17. 改进的PrefixSpan算法在Web挖掘中的应用](#)

[18. Web结构挖掘及HITS算法分析](#)

[19. Web结构挖掘与其基于超链接结构的算法](#)

[20. 基于Web结构挖掘算法的网站构建](#)

[21. Dijkstra算法的优化](#)

[22. Web结构挖掘](#)

[23. Dijkstra算法在企业成本控制的应用](#)

[24. 计算机行业挖掘结构性机会](#)

[25. Dijkstra算法在物流网络设计中的应用](#)

[26. 结构挖掘中web有向图模型的改进算法](#)

[27. Web结构挖掘研究](#)

[28. 基于Web页面链接结构的挖掘算法](#)

[29. XML的DOM树结构在WEB挖掘中的应用](#)

[30. 应用Web结构挖掘的PageRank算法的改进研究](#)

[31. 二进制挖掘算法在Web使用挖掘中的应用](#)

[32. Web结构挖掘中HITS算法的改进](#)

[33. 模式恢复算法在Web使用挖掘中的应用](#)

[34. 云计算在Web结构挖掘算法中的运用研究](#)

[35. Dijkstra算法在配电网抢修中的应用](#)

[36. Dijkstra算法在人群疏散上的应用](#)

[37. Dijkstra算法在路由选择中的应用](#)

[38. Dijkstra算法在最优投资策略问题中的应用](#)

[39. 基于Web结构挖掘算法的网站构建](#)

[40. 基于Web结构挖掘的HITS算法分析及改进](#)

[41. Dijkstra矩阵算法](#)

[42. Web结构挖掘研究](#)

[43. 云计算在Web结构挖掘算法中的运用研究](#)

[44. 数据聚类算法在web数据挖掘中的应用](#)

[45. Dijkstra算法在最短旅游路径中的应用](#)

[46. Dijkstra算法在部队快速行进中的应用](#)

[47. 基于PageRank和HITS的Web结构挖掘算法研究](#)

[48. 基于GIS的Dijkstra算法在运输系统的应用](#)

[49. 改进Dijkstra算法在PGIS中的应用](#)

[50. 基于Web结构挖掘的HITS算法研究](#)