

互联网数据挖掘作业报告——PageRank

学号：1400012783 姓名：王君珊

本次作业用 C++ 实现了 PageRank 算法, 并基于 ANN 论文数据集 (2013 年) 计算 Paper、Author 和 Venue 的 PageRank 值。

PageRank 是一种根据网页之间相互的超链接计算的技术, 作为网页排名的重要因素。

PageRank 的迭代计算公式如下： $\pi = \alpha P^T + (1 - \alpha) \frac{1}{n} e$, $e = (1, 1, \dots, 1)^T$

一、程序介绍：

本次作业有三个程序, 分别计算了 Paper、Author 和 Venue 的 PageRank 值。其中, Author 和 Venue 的程序结构相似, 而 Paper 的结构与前两者稍有不同。

1. PaperPageRank：

程序：PaperPageRank 源.cpp

输出文件：paperPagerank.txt

Paper 的 PageRank 值计算相对麻烦。首先, 由于 paper 的 ID 是字符串的形式, 且位数大, 故不能直接将其作为索引, 而是运用类似 Hash 的方法, 为其重新映射一个新的 ID。同时由于 paper 数量较大, 直接开辟 $N \times N$ 的 p 矩阵对内存要求较高, 故这里运用了类似稀疏矩阵的方法, 每个结点为一个结构体, 存储了指向它的结点及其 p 值。在每轮迭代计算 PageRank 值时, 遍历每个结点, 根据其所有父节点的 PageRank 值和相对应的 p 值, 求和得到新的 PageRank 值。

计算 Paper 的 PageRank 值用到了三个文件：paper_outcites.txt, acl.txt, paper_ids.txt。paper_outcites.txt 用来得到 p 矩阵的值, acl.txt 得到 paper 之间的引用关系, paper_ids.txt 得到 paper 的名字。

需要说明的是, 一共有 21212 篇 paper, 但是存在引用关系的只有 18164 篇, 也就是说有部分 paper 既没有引用也没有被引用, 程序仍然计算了它们的 PageRank 值。

2. AuthorPageRank：

程序：AuthorPageRank 源.cpp

输出文件：authorPagerank.txt

Author 的 PageRank 程序中的数据结构相对直观, 用一个 $N \times N$ 的矩阵直接存储 p 值, 每轮迭代直接运用文章开头的公式, 计算新的 PageRank 值。

计算 author 的 PageRank 值用到了两个文件：author_citation_network.txt, author_ids.txt。author_citation_network.txt 得到 author 之间的引用关系, author_ids.txt 得到 author 的名字。

需要注意的是, author_ids.txt 文件中的 ID 并不是连续的, 中间有部分 ID 号是被跳过的, 若直接以此 ID 作为索引会影响 PageRank 的计算, 故程序为它们重新安排序号。

3. VenuePageRank：

程序：VenuePageRank 源.cpp

输出文件：venuePagerank.txt

Venue 的 PageRank 程序结构与 author 的类似, 唯一不同的是 venue 之间的引用关系需要通过 paper 的引用关系获得。所以需要进行数据预处理, 即得到每篇 paper 所属的 venue。数据预处理后, 用类似 Paper 的输入方法读入 paper 的引用关系, 转换成 venue 之间的引用关系, 用类似 author 的迭代方法计算。

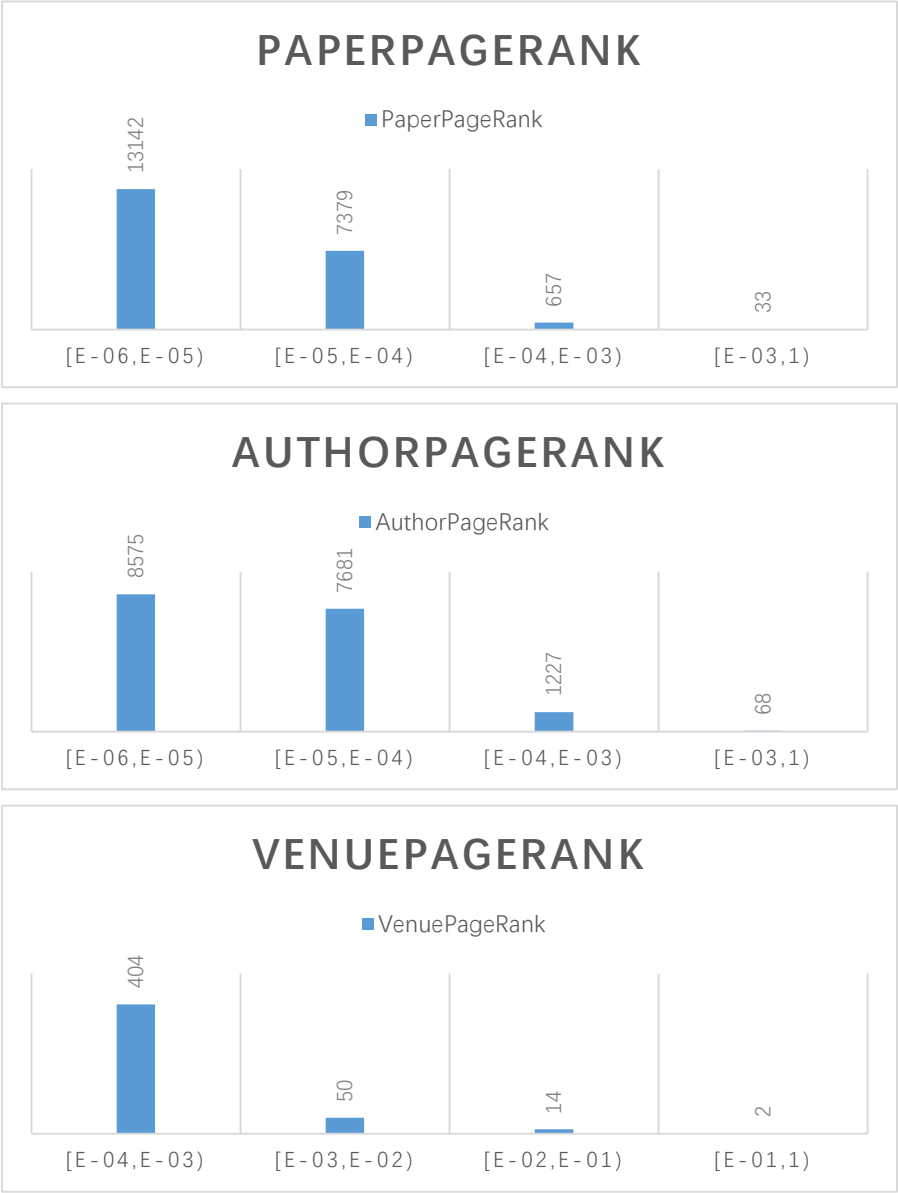
计算 venue 的 PageRank 值用到了两个文件：acl.txt, acl-metadata.txt。acl.txt 得到 paper 之间的引用关系, 进而得到 venue 的引用关系, acl-metadata.txt 得到 paper 和 venue 的对

应关系和 venue 的名字。

需要注意的是，对于两个 venue 之间多篇 paper 引用的情况，处理方法是相当于两个 venue 之间有多条边，p 值相对应进行加权，即引用越多，p 值会越大。

二、数据汇报与分析分析：

排序后的 PaperPageRank、AuthorPageRank 和 VenuePageRank 存在相应的 txt 文件中，各自 PageRank 的分布情况如下表：



可以看出，无论是 paper、author 还是 venue，数量都会随着 PageRank 值变大而减少，也就是说大部分的 PageRank 值都在较低水平，只有极少量的 Paper、Author 或者 Venue 有较高的 PageRank 值，即有较高的引用量，处于顶尖水平。

三、总结

本次作业的难点我认为有主要两点：一是对文本文件的处理，即如何从大量格式不规范的数据中获得有用信息，并转换为可以计算的数据；二是对于大量的数据，如何在保证计算正确的前提下，降低时间和空间的复杂度。