

互联网数据挖掘第一次作业

——pagerank 计算

朱政烨 1700017760

数据处理

- 首先从 wikidump 下载 60+GB 的 xml 文件
- 解析 xml 的方法有很多, 如 DOM、SAX、ElementTree 等, 但是我这次采用了 ElementTree, 它与 DOM 相比速度更快, 与 SAX 相比同样可以用 iterparse 函数, 不会一次性在内存中读入整个文档。这对于本次大文件读取是关键。
- 具体来说, 本次利用 python 的 lxml 库, 使用 etree 的 iterparse 方法, 逐个读取文件中的每个 element 并获取其标签。
- 经过研究数据结构, 发现对于本次计算有用的标签有下列几个:
 - <page>代表一个词条的整个页面
 - <title>是词条的名称
 - <id>是词条的序号 (本次不采用, 而是直接用 title 代表一个词条)
 - <text>是词条正文, 我们要求的外链就在此中, 并且处于双中括号[]中
- 我们采用循环遍历每个 tag, 每到一个 page 处, 就记录下 title, 并从 text 中用正则表达式提取出外链, 最后都存储下来, 存储格式为字典{"title1":["外链 1","外链 2",...], "title2":["外链 1","外链 2",...], ...}, 其实就是我们熟悉的邻接表结构。

算法描述

- pagerank 算法其实很简单, 在我们上一步能够得出转移矩阵 P 后, 直接用 P^T 与 pagerank 值向量 π 相乘, 反复多次后就能收敛。(马尔科夫过程)

— Transition Probability Matrix $P = \{p_{ij}\}$

$$p_{ij} = \begin{cases} \frac{M(i,j)}{\sum_{k \in \text{outlink}[v_i]} M(i,k)}, & \text{outlink}[v_i] \neq \emptyset \\ M(i,j) = 0, & \text{otherwise} \end{cases}$$

$$P = \begin{bmatrix} 0 & 1/3 & 1/3 & 0 & 0 & 1/3 & 0 & 0 \\ 0 & 0 & 0 & 1/2 & 1/2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1/2 & 0 & 1/2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 1/4 & 0 & 0 & 1/4 & 0 & 0 & 1/4 & 1/4 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1/2 & 0 & 0 & 1/2 & 0 \end{bmatrix}$$

$$\pi(v_i) = \sum_{v_j \in \text{inlink}[v_i]} \frac{\pi(v_j)}{|\text{outlink}[v_j]|} \quad \longrightarrow \quad \pi = P^T \pi$$

- 我们这里主要处理这个矩阵向量乘法。直接在 python 中乘是不现实的: 首先这个矩阵是 $10^6 \times 10^6$ 的, 内存不够肯定不能直接运行; 而且它也是稀疏的, 之前已经用邻接表结构存储, 完全没有必要用这个大矩阵。

- 于是我回想到之前学的并行与分布式计算，曾经提到了矩阵向量乘法不同的分解方法：
(按行、按列、按块)

$$\begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n \\ \vdots \\ a_{n1}x_1 + a_{n2}x_2 + \cdots + a_{nn}x_n \end{bmatrix}$$

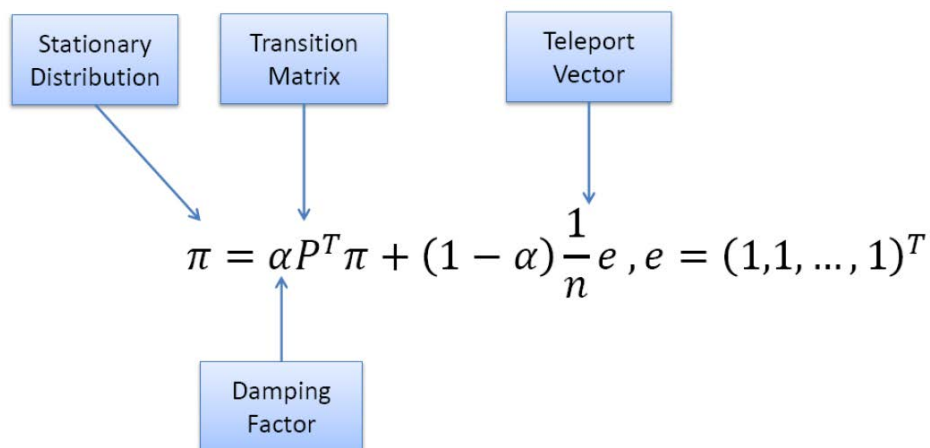
传统的按行分解

$$\begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n \\ \vdots \\ a_{n1}x_1 + a_{n2}x_2 + \cdots + a_{nn}x_n \end{bmatrix}$$

按列分解

那我们这里显然可以采用第二种方式。矩阵的第一列，转置前就是第一行，这一行我们很熟悉——它就是我们已经存在邻接表里的 x_1 网页相应的出链，把他与 x_1 相乘，就能得到 x_1 对所有页面的贡献，也就是乘后向量的第一列，然后 x_2 、 x_3 ...都这样依次相加，就得到了全部页面对全部页面的贡献，数学上也就是图中的按列相乘，列与列相加的最终结果。

- 实际操作后发现，我们采样的数据量很小，存在排序泄露和排序沉入问题，导致收敛速度很慢而且 pr 值不合理（迭代 20、50、100 次的结果相差很大，而且少数值过大，以及大量的 0 值）。于是我加入了 RWR 的改进，在随机游走过程中重新浏览新网页，也就是在原来的迭代中加入阻尼因数 $\alpha=0.85$ 。事实证明这条改动是大有用处的。



结果分析

- 收敛情况

分别设置迭代次数为 1、10、20、50，输出到文件中，观察数值变化及收敛情况。

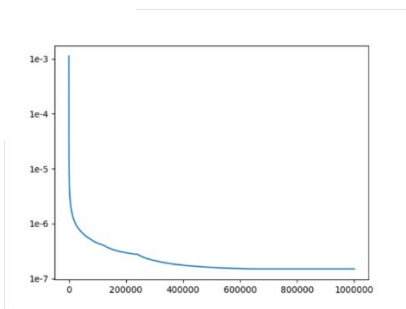
下图依次是 1、10、20、50 次迭代后第一页的词条及结果，可以发现，10 次后已经趋于稳定，而 50 次的结果已经和 20 次相差无几（pr 值小数点九位前都没有任何变化，100 万个词条的排名也完全一样），至此已经可以得出收敛的结论。

1 United States Census Bureau 0.0021208038152328614	1 Race and ethnicity in the United States Census 0.0011158697870010288
2 Marriage 0.0007644432649006158	2 United States Census Bureau 0.0008670301294969095
3 White (U.S. Census) 0.0007405483830640203	3 Marriage 0.00027190410511977255
4 List of sovereign states 0.0006628857399271051	4 White (U.S. Census) 0.0002626858749363128
5 Hispanic (U.S. Census) 0.0006055275423950848	5 List of sovereign states 0.0002187586532843094
6 Latino (U.S. Census) 0.0006042817372393482	6 Hispanic (U.S. Census) 0.00021713113773311798
7 Geographic Names Information System 0.0005412706729629829	7 Latino (U.S. Census) 0.0002166041647986756
8 Native American (U.S. Census) 0.000513971944494449	8 Geographic Names Information System 0.0001912515740766653
9 World War II 0.0004739688027510905	9 Native American (U.S. Census) 0.00018970803682307768
10 African American (U.S. Census) 0.00047234708070642545	10 African American (U.S. Census) 0.00017575253509664993
11 Japan 0.00046769421303918547	11 The New York Times 0.00017287527045143058
12 Asian (U.S. Census) 0.00045456427551658895	12 World War II 0.00017230183669815233
13 United Kingdom 0.0004375972846591556	13 Asian (U.S. Census) 0.00017131339644666918
14 Germany 0.0004186649339423988	14 United Kingdom 0.0001458862918105785
15 England 0.00040131399472738107	15 Japan 0.00014026075473392866
16 New York City 0.00034446312056895866	16 New York City 0.0001274536718690047
17 The New York Times 0.0003227239605476139	17 Germany 0.00012106529313986912
18 London 0.0002974570664889391	18 England 0.00010655651573307543
19 India 0.0002580606059830049	19 London 5.408977226777743e-05
20 California 0.0002407683086071506	20 Pennsylvania 8.815025459483422e-05
21 Italy 0.00023863072609571214	21 California 8.578857306741958e-05
22 Brazil 0.000232539509105412	22 Republican Party (United States) 8.296767714164906e-05
23 Pennsylvania 0.00022620629859764946	23 India 7.991908361766908e-05
24 Netherlands 0.00021244881590426425	24 Italy 7.622330506146213e-05
25 Minnesota 0.00020359044668125793	25 Soviet Union 7.403000546328376e-05
26 Poland 0.0002015261507353246	26 Washington, D.C. 7.344804001070421e-05
27 China 0.0001906016096598456	27 Minnesota 7.292827564979789e-05
28 Ontario 0.000189244171169571	28 China 6.9815329056142e-05
29 Paris 0.0001863733180210424	29 Catholic Church 6.752408752495312e-05
30 Communes of France 0.00018134595601412943	30 Paris 6.629506310629072e-05
31 Russia 0.00017187555687787358	31 Latin 6.600537921711388e-05
32 Scotland 0.00017062860464462875	32 Netherlands 6.517103737321558e-05
33 Soviet Union 0.00016794200343129466	33 Russia 6.366781196163792e-05
34 Spain 0.00016741362625080496	34 Europe 6.259786986564783e-05
35 Republican Party (United States) 0.00016100579796920013	35 Poland 6.036905802982225e-05
36 List of Enochian angels 0.0001566066666666634	36 Texas 6.0110769916304016e-05
37 Texas 0.000153729533198113	37 Pacific Islander (U.S. Census) 5.95677007426111e-05

1 Race and ethnicity in the United States Census 0.0011149293506275526	1 Race and ethnicity in the United States Census 0.0011149291388117756
2 United States Census Bureau 0.0008665013205450971	2 United States Census Bureau 0.0008665012009306858
3 Marriage 0.000271904765699448	3 Marriage 0.0002719047503062613
4 White (U.S. Census) 0.0002626099583058004	4 White (U.S. Census) 0.0002626099583058004
5 List of sovereign states 0.00021866260449354947	5 List of sovereign states 0.00021866277556452636
6 Hispanic (U.S. Census) 0.000217119709755478	6 Hispanic (U.S. Census) 0.0002171196934931443
7 Latino (U.S. Census) 0.00021653302020666922	7 Latino (U.S. Census) 0.0002165330037319481
8 Geographic Names Information System 0.00019118367420095506	8 Geographic Names Information System 0.0001911836584876276
9 Native American (U.S. Census) 0.0001896407340065376	9 Native American (U.S. Census) 0.0001896407184163799
10 African American (U.S. Census) 0.00017560707164549566	10 African American (U.S. Census) 0.00017560705640520943
11 The New York Times 0.00017234046242221046	11 The New York Times 0.0001723403266510781
12 World War II 0.00017195567337449638	12 World War II 0.00017195556761429612
13 Asian (U.S. Census) 0.00017124732280174533	13 Asian (U.S. Census) 0.00017124730754909603
14 United Kingdom 0.0001453254650851864	14 United Kingdom 0.0001453253865608537
15 Japan 0.0001399683352127016	15 Japan 0.0001399665553980116
16 New York City 0.00012716982281384167	16 New York City 0.00012716975042072174
17 Germany 0.00012086870524355158	17 Germany 0.00012086864374654629
18 England 0.0001064080929566037	18 England 0.00010640886331644641
19 London 9.469129838356362e-05	19 London 9.469123767060948e-05
20 Pennsylvania 8.804585779252049e-05	20 Pennsylvania 8.804593243104637e-05
21 California 8.564571297944051e-05	21 California 8.5645790918223e-05
22 Republican Party (United States) 8.276446331890432e-05	22 Republican Party (United States) 8.276441611715596e-05
23 India 7.98005707773549e-05	23 India 7.980053651361685e-05
24 Italy 7.60561710306101e-05	24 Italy 7.605611696832542e-05
25 Soviet Union 7.384758935997144e-05	25 Soviet Union 7.38475362915656e-05
26 Washington, D.C. 7.31479158234333e-05	26 Washington, D.C. 7.314784421913962e-05
27 Minnesota 7.287845351566879e-05	27 Minnesota 7.287844185074531e-05
28 China 6.968454171417176e-05	28 China 6.96845044845325e-05
29 Catholic Church 6.720450359338822e-05	29 Catholic Church 6.720440279760835e-05
30 Paris 6.61589960931748e-05	30 Paris 6.615865660420637e-05
31 Latin 6.57910886054914e-05	31 Latin 6.5791018737684e-05
32 Netherlands 6.503862891285966e-05	32 Netherlands 6.503858754096144e-05
33 Russia 6.352549442788358e-05	33 Russia 6.352545198901926e-05
34 Europe 6.247446125783996e-05	34 Europe 6.247442591019842e-05
35 Poland 6.026344934429346e-05	35 Poland 6.026341717610626e-05
36 Texas 6.002401192773314e-05	36 Texas 6.0023991515087514e-05
37 Pacific Islander (U.S. Census) 5.953103630555582e-05	37 Pacific Islander (U.S. Census) 5.953103630555582e-05

- 结果分布

通过 matplotlib 绘图，得到数据的分布图。可以发现，pr 值虽然从 1e-7 到 1e-3 都有，但是主要还是集中在 1e-7~1e-6 之间，高于 1e-6 的仅有 23130 个词条，高于 1e-5 的仅有 589 条。



- 热门词条

我选取了 100 万中排名前 100 个词条，翻译成中文，可以看出美国人口普查相关内容占很大的比重，而且其中人种中白人排第一。除此之外，英国是排名最高的国家，纽约是排名最高的城市。当然，我们只在 wiki 里选了 100 万个词条，得到的这些信息也只能代表前 100 万词条的趋势罢了。

1	美国人口普查中的种族和种族
2	美国人口普查局
3	结婚
4	白人（美国人口普查）
5	主权国家名单
6	西班牙裔（美国人口普查）
7	拉丁美洲人（美国人口普查）
8	地名信息系统
9	美洲原住民（美国人口普查）
10	非裔美国人（美国人口普查）
11	纽约时报
12	第二次世界大战
13	亚洲人（美国人口普查）
14	英国
15	日本
16	纽约城
17	德国
18	英格兰
19	伦敦
20	宾夕法尼亚州
21	加利福尼亚州
22	共和党（美国）
23	印度
24	意大利
25	苏联
26	华盛顿特区
27	明尼苏达州
28	中国
29	天主教堂
30	巴黎



总结反思

通过这次作业，我掌握了初步的数据挖掘知识，包括 xml 文件的解析、pagerank 算法的简单实现，编程能力也有了进一步的提升。这次作业还能有很多提升的地方，比如算法的优化、并行化，数据的优化等等，希望下次时间充裕可以更上一层楼。

2019.10.8