

# 目录

1. 实验说明.....	2
1.1 影响最大化介绍.....	2
1.2 独立级联模型.....	2
1.3 线性阈值模型.....	2
2. 实验方法.....	2
2.1 寻找自洽排名.....	3
2.2 加上贪婪算法的改进.....	4
3. 实验过程.....	4
3.1 读取文件.....	4
3.2 IMRank 算法实现.....	5
3.3 模拟爆发.....	6
3.4 算法改进.....	6
3.5 主函数.....	8
4. 实验结果.....	8
4.1 中间变量.....	8
4.2 最终结果.....	9

# 1. 实验说明

## 1. 影响最大化介绍

社交网络可以抽象的用图  $G(V, E, P)$  来描述， $V$  是结点集， $E$  是边集， $P$  是边的概率集。一个用户就是一个节点  $v$ ，用户与用户之间的关系就是边  $e$ ，每条边都有一条概率  $p$ ，感染会根据边上的概率进行传播。

在此描述下，影响最大化问题可以分为两类：

1. 给定节点数，选择出  $k$  个节点作为种子集，使得种子集能影响的节点数尽可能多。
2. 给定所要求产生的影响力，找到满足条件的最小节点集合

在本实验中，研究的问题为第一类，在实验种，分别选取 5、10、15、20 个节点，使得选取的节点，在给定的节点个数中，可以影响尽可能多的节点数。

## 1.2 独立级联模型

独立级联模型（Independent Cascade Model）简称 IC 模型，是一种概率传播模型，基本原理描述如下：

在社交网络  $G=(V, E)$  中，点集  $V$  中的节点具有两种状态，一种是已感染状态，一种是未感染状态。每一个处于已感染状态的结点，都有按边权的概率感染每个与之相连的结点的机会。每个节点只有一次被相邻节点尝试感染的机会。在初始状态时，仅有种子节点集合中的点处于已感染状态，这些感染的点按照概率尝试感染与之相邻的点，直至集合中所有已感染的点，都尝试感染过与之相邻的所有点，并且无新增的已感染点。

## 1.3 线性阈值模型

线性阈值模型与独立级联模型的区别表现在，在线行阈值模型中，每个节点  $v$  受到相邻已感染节点的影响的叠加，当超过某个阈值时，该节点被感染。

# 2. 实验方法

## 2.1 寻找自洽排名

本次实验中采用的是独立级联模型。

实验中首先采用迭代的方式,通过寻找自适应的排名来计算给定图的影响最大化的点集<sup>[1]</sup>。具体算法如算法 1 和算法 2。

### 算法 1:

过程: 更新边界影响力 Calculate  $Mr(r)$

输入:  $Mr(r)$ , 图  $G(V, E)$ , 其中边的权重  $p(v_{rj}, v_{ri})$  表示从  $v_{rj}$  直接感染  $v_{ri}$  的概率

输出: 更新之后的  $Mr(r)$

1. for  $i = 1$  to  $n$  do
2.      $Mr(v_{ri}) \leftarrow 1$
3. end for
4. for  $i = 2$  to  $n$  do
5.     for  $j = 1$  to  $i$  do
6.          $Mr(v_{rj}) \leftarrow Mr(v_{rj}) + p(v_{rj}, v_{ri}) * Mr(v_{ri})$
7.          $Mr(v_{ri}) \leftarrow (1 - p(v_{rj}, v_{ri})) * Mr(v_{ri})$
8.     end for
9. end for
10. output  $Mr$

### 算法 2:

过程: 输出边界影响序列  $IMRank(r)$

输入: 图  $G(V, E)$ , 初始排名序列  $r$

输出: 边界影响力排名  $r$

1.  $r^{(0)} \leftarrow r$
2.  $t \leftarrow 0$
3. 设置初始时所有节点的边界影响力均为 1
4. 设置初始时所有节点的排名  $r^{(0)}$  (全为 1)
5. repeat
6.      $t \leftarrow t + 1$

7. 调用更新边界影响力过程，将边界影响力更新为  $Mr(t)$
8. 通过对边界影响力进行排序，生成新的边界影响力的排名  $r^{(t)}$
9. until  $r^{(t)} = r^{(t-1)}$
10. output  $r^{(t)}$

其中采用的符号表如图 2.1 所示

Notation	Description
$v_i$	a node with index $i$
$r_i$	the index of node with rank $i$ with respect to a given ranking $r$
$S = \{v_1, v_2, \dots, v_n\}$	a set of nodes
$I(S)$	expected number of nodes eventually activated by set $S$
$M(v S)$	marginal influence spread by adding node $v$ into a seed set $S$
$M_r(v_{r_i})$	ranking-based marginal influence, short for $M(v \{v_{r_1}, v_{r_2}, \dots, v_{r_{i-1}}\})$
$p(v_i \{v_1, v_2, \dots, v_{i-1}\})$	probability that $v_i$ is activated given that a collection of nodes $\{v_1, v_2, \dots, v_{i-1}\}$ are already activated
$\eta_r(v_i, v_j)$	influence score that node $v_i$ sends to node $v_j$ with respect to a given ranking $r$
$d(v_j, v_i)$	a simple path starting from $v_j$ and ending at $v_i$ , i.e., $\{w_1 = v_j, w_2, \dots, w_n = v_i\}$
$d_r(v_j, v_i)$	influence path, which is a simple path where $v_j$ is the only node ranked higher than $v_i$ on the path
$\rho_r(v_i, v_j)$	probability that $v_i$ is activated by $v_j$ through any influence path, with respect to a given ranking $r$
$l$	maximal length of all influence paths to account into

图 2.1 符号表

## 2.2 加上贪婪算法的改进

在实验过程中，初始排名不同时，计算所得到的排名可能不一致。在实验中采用贪婪算法解决这个问题。实验中采用两个完全相反的排名初始序列，根据者两个不同的初始排名，可能得到两个不完全相同的 top k 点序列。对于这两个序列中相同的前一部分，认为对给定数据集而言，这重合的一部分一定在 top k 点序列中。对于后面不相同的部分，则考虑是因为数据中边的权重设置，使得该算法不能有效区分这几个点的边界影响力。因此对于这些点，采用模拟爆发一定次数，来计算该节点的影响力，然后将较大影响力的点，加入 top k 点序列，直至序列长度达到 k。具体操作为：

有两个根据不同初始排名得到的两个不同的序列 t1、t2。采用两个指针 i、j 分别扫描 t1、t2。若指向相同的结点，则说明此结点排序正确，则将 i 和 j 都指向列表中下一个元素。若指向不同的结点，则分别采用模拟爆发的方法来计算将元素加入已经正确排序的结点的序列的影响力，然后选择较大的影响力对应的节点，认为该节点更适合加入 top k 点序列。

## 3. 实验过程

### 3.1 读取文件

首先从磁盘中读取数据文件，将其用邻接矩阵的形式保存下来

```
1. def load_file(file_name="DUNF with Weights.txt"): # 读取文件，文件形式为三列，
    分别为父节点、子节点、权重，返回图的邻接矩阵
2.     with open(file_name, 'r') as fp:
3.         lines = fp.readlines()
4.         num = int(max(max([float(_) for _ in line.split()]) for line in lines)) # 获取节点的个数
5.         data = np.zeros((num+1, num+1), np.float32)
6.         for line in lines:
7.             data[int(line.split()[0])][int(line.split()[1])] = float(line.split()[2])
8.     return data
```

### 3.2 IMRank 算法实现

采用 last to first 策略，算法 1 的实现代码如下

```
1. def last_to_first(mi_rank, data): # 输入数据和现有排名，计算下一次迭代产生的边界影响力
2.     mi_list = [1 for _ in range(len(mi_rank))] # 存放每个结点对应的  $Mr(r)$ ，表示该点的边界影响力
3.     for i in range(len(mi_rank)-1, -1, -1): # 取出排名为 i 的结点
4.         for j in range(i):
5.             mi_list[mi_rank[j]] = mi_list[mi_rank[j]] + data[mi_rank[j]][mi_rank[i]] * mi_list[mi_rank[i]]
6.             mi_list[mi_rank[i]] = (1 - data[mi_rank[j]][mi_rank[i]]) * mi_list[mi_rank[i]]
7.     return mi_list
```

算法 2 的实现代码如下

```
1. def im_rank(data, init_rank, top_k=20, iter_num=40): # 返回的边界影响力是全部点的影响力
2.
3.     mi_rank_1 = [_ for _ in init_rank] # 第零次迭代产生的排名
4.     iter_count = 0 # 迭代次数
5.     equal = False
6.     while not equal and iter_count <= iter_num: # 终止迭代条件（满足其一）：1. 迭代不导致结果变动 2. 达到一定迭代次数
7.         tmp = last_to_first(mi_rank_1, data)
8.         mi_rank_2 = np.argsort(np.array(tmp))[:-1].tolist()
```

```

9.         equal = (np.array(mi_rank_1) == np.array(mi_rank_2)).all() # 如果两
           次迭代得到的边界影响最大的 k 个点相同，则停止迭代
10.         mi_rank_1 = [_ for _ in mi_rank_2]
11.         iter_count = iter_count + 1
12.     return mi_rank_2[:top_k:]

```

### 3.3 模拟爆发

采用 IC 模型模拟爆发，返回感染的字典

```

1. def ic_model(data, seed_list, total=100): # 输入带权图、爆发种子结点、爆发次
     数，返回键为感染点、值为感染几率的字典
2.     active_dict = {}
3.     for break_num in range(total):
4.         active_set = np.zeros(len(data), np.int) # 标记点的感染状态，为 0 表示
           未感染，为 1 表示感染
5.         child_queue = queue.Queue() # 待考察点
6.         for _ in seed_list:
7.             child_queue.put(_)
8.         while not child_queue.empty(): # 当队列不为空时，即存在已感染的点的孩子
           结点尚未考察
9.             i = child_queue.get() # 本次考察的点
10.            for j in range(len(data)): # 对于该点的每个孩子
11.                if j != i and data[i][j] != 0: # 不为自己，也不为零
12.                    if active_set[i] != 1 and data[i][j] > random.random():
                       # 在此次传染中，由未感染变为感染
13.                        if j not in active_dict:
14.                            active_dict[j] = 1
15.                        else:
16.                            active_dict[j] += 1
17.                            active_set[i] = 1 # 设置为被感染的状态
18.                            child_queue.put(j) # 加入待考察队列
19.        for _ in active_dict:
20.            active_dict[_] = active_dict[_] / total
21.    return active_dict

```

### 3.4 算法改进

采用贪婪算法，对于两个不同的序列，不断的从这两个序列中挑选影响最大的结点，直到已经挑选了 k 个点时停止。

```

1. def greedy_mi_influence(order1, order2, k, data): # 传入两次迭代产生的结果，返回模拟爆发之后的序列
2.     result = []
3.     i = 0
4.     j = 0
5.     while len(result) < k:
6.         if order1[i] in result:
7.             i = i + 1
8.             continue
9.         if order2[j] in result:
10.            j = j + 1
11.            continue
12.         if order1[i] == order2[j]: # 如果两处排名相同，则说明排名没有波动，则加入序列
13.             result.append(order1[i])
14.             i = i + 1
15.             j = j + 1
16.         else:
17.             tmp1 = sum(ic_model(data, result+order1[i:i+1:1]).values()) # 计算 order1[i] 的边界影响力
18.             tmp2 = sum(ic_model(data, result+order2[j:j+1:1]).values()) # 计算 order2[j] 的边界影响力
19.             if tmp1 > tmp2:
20.                 result.append(order1[i])
21.                 i = i + 1
22.             else:
23.                 result.append(order2[j])
24.                 j = j + 1
25.     return result

```

结合两种方法，计算排名

```

1. def combine_rank(data, top_k=20, iter_num=30): # 综合采用两种方法计算影响力排名
2.     init_rank_1 = [len(data) - i - 1 for i in range(len(data))] # 生成的初始排名 1
3.     init_rank_2 = [i for i in range(len(data))] # 生成的初始排名 2
4.     mi_rank_1 = im_rank(data, init_rank_1, top_k, iter_num)
5.     mi_rank_2 = im_rank(data, init_rank_2, top_k, iter_num)
6.     mi_order = greedy_mi_influence(mi_rank_1, mi_rank_2, top_k, data) # 迭代可能未收敛，采用贪婪算法再次处理
7.     return mi_order

```

调用上面函数，计算最大影响力的函数

```

1.
2. def max_influence(data, top_k=5, iter_num=30, break_num=10): # 返回影响力最大的 k 个点, 和这 k 个点在模拟爆发时的平均感染点数
3.     begin_time = time() # 开始计算边界影响力的时间
4.     order = combine_rank(data, top_k, iter_num)
5.     end_mi_cal = time() # 边界影响力计算完成的时间
6.     duration1 = end_mi_cal - begin_time
7.     break_result = ic_model(data, order, break_num) # 传入用于爆发的种子时, 应该从 order 列表中选取前 k 个, k 的选取为计算边界影响最大的的个数
8.     duration2 = time() - end_mi_cal # 计算模拟爆发所用时长
9.     print("计算影响力最大的%d 个点用时: %fs\n 计算这%d 点模拟爆发用时: %fs" % (top_k, duration1, top_k, duration2))
10.    active_count = 0
11.    for _ in break_result:
12.        active_count += break_result[_]
13.    return order, break_result

```

### 3.5 主函数

运行主函数。计算的结果不方便

```

1. if __name__ == "__main__":
2.     Data = load_file("DUNF with Weights.txt")
3.     seedSetList = [] # 存放每次计算得到影响力最大的 k 个点
4.     activeSetList = [] # 存放感染的
5.     activeCountList = [] # 存放平均感染的点的个数
6.     iterNum = 30 # 计算边界影响力时的迭代次数上界
7.     breakNum = 100 # 模拟爆发点的次数
8.     for topK in list(range(5, 25, 5)):
9.         seedSet, activeSet = max_influence(Data, topK, iterNum, breakNum)
10.        seedSetList.append(seedSet)
11.        activeSetList.append(activeSet)
12.        activeCountList.append(sum(activeSet.values()))
13.
14.    print("end")

```

## 4. 实验结果

### 4.1 中间变量

在调试模式下, 查看运行的中间变量, 调用贪婪算法之后得到的结果如图 4.1 到图 4.4 所示。其中 init\_rank\_1 和 init\_rank\_2 表示两个不同的初始排



名，mi\_list\_1 和 mi\_list\_2 分别对应这两个初始排名对应的 IMRank 算法得到的结果，mi\_order 对应调用贪婪算法之后得到的排名。

```
▶ In [5]: mi_order = {list: 5} [746, 592, 732, 470, 90]
▶ In [6]: mi_rank_1 = {list: 5} [746, 592, 732, 470, 559]
▶ In [7]: mi_rank_2 = {list: 5} [90, 245, 67, 131, 136]
In [8]: top_k = {int} 5
```

图 4.1 Top 5 个结点

```
▶ In [9]: mi_order = {list: 10} [746, 592, 732, 90, 245, 470, 67, 131, 136, 301]
▶ In [10]: mi_rank_1 = {list: 10} [746, 592, 732, 470, 559, 677, 658, 546, 581, 740]
▶ In [11]: mi_rank_2 = {list: 10} [90, 245, 67, 131, 136, 301, 292, 107, 47, 88]
In [12]: top_k = {int} 10
```

图 4.2 Top 10 个结点

```
▶ In [13]: mi_order = {list: 15} [746, 90, 245, 592, 732, 67, 131, 136, 470, 559, 677, 301, 658, 292, 107]
▶ In [14]: mi_rank_1 = {list: 15} [746, 592, 732, 470, 559, 677, 658, 546, 581, 740, 604, 626, 698, 434, 629]
▶ In [15]: mi_rank_2 = {list: 15} [90, 245, 67, 131, 136, 301, 292, 107, 47, 88, 28, 161, 295, 65, 111]
In [16]: top_k = {int} 15
```

图 4.3 Top 15 个结点

```
▶ In [17]: data = (ndarray: (751, 751)) [[0. 0. 0. ... 0. 0. 0.], [0. 0. 0. ... 0. 0. 0.], [0. 0. 0. ... 0. 0. 0.], ..., [0. 0. 0. ... 0. 0. 0.], [0. 0. 0. ... 0. 0. 0.], [0. 0. 0. ... 0. 0. 0.]] ...View as Array
▶ In [18]: init_rank_1 = {list: 751} [750, 749, 748, 747, 746, 745, 744, 743, 742, 741, 740, 739, 738, 737, 736, 735, 734, 733, 732, 731, 730, 729, 728, 727, 726, 725, 724, 723, 722, 721, 720, 719, 718, 717, 716, 715, 714, 713, 712, 711, 710, 709, 708, 707, 706, 705, 704, 703, 702, 701, 700, 699, 698, 697, 696, 695, 694, 693, 692, 691, 690, 689, 688, 687, 686, 685, 684, 683, 682, 681, 680, 679, 678, 677, 676, 675, 674, 673, 672, 671, 670, 669, 668, 667, 666, 665, 664, 663, 662, 661, 660, 659, 658, 657, 656, 655, 654, 653, 652, 651, 650, 649, 648, 647, 646, 645, 644, 643, 642, 641, 640, 639, 638, 637, 636, 635, 634, 633, 632, 631, 630, 629, 628, 627, 626, 625, 624, 623, 622, 621, 620, 619, 618, 617, 616, 615, 614, 613, 612, 611, 610, 609, 608, 607, 606, 605, 604, 603, 602, 601, 600, 599, 598, 597, 596, 595, 594, 593, 592, 591, 590, 589, 588, 587, 586, 585, 584, 583, 582, 581, 580, 579, 578, 577, 576, 575, 574, 573, 572, 571, 570, 569, 568, 567, 566, 565, 564, 563, 562, 561, 560, 559, 558, 557, 556, 555, 554, 553, 552, 551, 550, 549, 548, 547, 546, 545, 544, 543, 542, 541, 540, 539, 538, 537, 536, 535, 534, 533, 532, 531, 530, 529, 528, 527, 526, 525, 524, 523, 522, 521, 520, 519, 518, 517, 516, 515, 514, 513, 512, 511, 510, 509, 508, 507, 506, 505, 504, 503, 502, 501, 500, 499, 498, 497, 496, 495, 494, 493, 492, 491, 490, 489, 488, 487, 486, 485, 484, 483, 482, 481, 480, 479, 478, 477, 476, 475, 474, 473, 472, 471, 470, 469, 468, 467, 466, 465, 464, 463, 462, 461, 460, 459, 458, 457, 456, 455, 454, 453, 452, 451, 450, 449, 448, 447, 446, 445, 444, 443, 442, 441, 440, 439, 438, 437, 436, 435, 434, 433, 432, 431, 430, 429, 428, 427, 426, 425, 424, 423, 422, 421, 420, 419, 418, 417, 416, 415, 414, 413, 412, 411, 410, 409, 408, 407, 406, 405, 404, 403, 402, 401, 400, 399, 398, 397, 396, 395, 394, 393, 392, 391, 390, 389, 388, 387, 386, 385, 384, 383, 382, 381, 380, 379, 378, 377, 376, 375, 374, 373, 372, 371, 370, 369, 368, 367, 366, 365, 364, 363, 362, 361, 360, 359, 358, 357, 356, 355, 354, 353, 352, 351, 350, 349, 348, 347, 346, 345, 344, 343, 342, 341, 340, 339, 338, 337, 336, 335, 334, 333, 332, 331, 330, 329, 328, 327, 326, 325, 324, 323, 322, 321, 320, 319, 318, 317, 316, 315, 314, 313, 312, 311, 310, 309, 308, 307, 306, 305, 304, 303, 302, 301, 300, 299, 298, 297, 296, 295, 294, 293, 292, 291, 290, 289, 288, 287, 286, 285, 284, 283, 282, 281, 280, 279, 278, 277, 276, 275, 274, 273, 272, 271, 270, 269, 268, 267, 266, 265, 264, 263, 262, 261, 260, 259, 258, 257, 256, 255, 254, 253, 252, 251, 250, 249, 248, 247, 246, 245, 244, 243, 242, 241, 240, 239, 238, 237, 236, 235, 234, 233, 232, 231, 230, 229, 228, 227, 226, 225, 224, 223, 222, 221, 220, 219, 218, 217, 216, 215, 214, 213, 212, 211, 210, 209, 208, 207, 206, 205, 204, 203, 202, 201, 200, 199, 198, 197, 196, 195, 194, 193, 192, 191, 190, 189, 188, 187, 186, 185, 184, 183, 182, 181, 180, 179, 178, 177, 176, 175, 174, 173, 172, 171, 170, 169, 168, 167, 166, 165, 164, 163, 162, 161, 160, 159, 158, 157, 156, 155, 154, 153, 152, 151, 150, 149, 148, 147, 146, 145, 144, 143, 142, 141, 140, 139, 138, 137, 136, 135, 134, 133, 132, 131, 130, 129, 128, 127, 126, 125, 124, 123, 122, 121, 120, 119, 118, 117, 116, 115, 114, 113, 112, 111, 110, 109, 108, 107, 106, 105, 104, 103, 102, 101, 100, 99, 98, 97, 96, 95, 94, 93, 92, 91, 90, 89, 88, 87, 86, 85, 84, 83, 82, 81, 80, 79, 78, 77, 76, 75, 74, 73, 72, 71, 70, 69, 68, 67, 66, 65, 64, 63, 62, 61, 60, 59, 58, 57, 56, 55, 54, 53, 52, 51, 50, 49, 48, 47, 46, 45, 44, 43, 42, 41, 40, 39, 38, 37, 36, 35, 34, 33, 32, 31, 30, 29, 28, 27, 26, 25, 24, 23, 22, 21, 20, 19, 18, 17, 16, 15, 14, 13, 12, 11, 10, 9, 8, 7, 6, 5, 4, 3, 2, 1, 0]
In [19]: iter_num = {int} 30
▶ In [20]: mi_order = {list: 20} [746, 592, 732, 90, 245, 67, 131, 136, 470, 559, 677, 658, 301, 546, 292, 107, 47, 581, 740, 88]
▶ In [21]: mi_rank_1 = {list: 20} [746, 592, 732, 470, 559, 677, 658, 546, 581, 740, 604, 626, 698, 434, 629, 688, 534, 596, 737, 693]
▶ In [22]: mi_rank_2 = {list: 20} [90, 245, 67, 131, 136, 301, 292, 107, 47, 88, 28, 161, 295, 65, 111, 24, 432, 336, 39, 254]
In [23]: top_k = {int} 20
```

图 4.4 Top 20 个结点

4.2 最终结果

分别计算 top 5、10、15、20 个点集时的感染情况，采用如下记录方式：初始感染的种子节点为计算所得的 top 节点。计算结果不包含初始被选中的节点

$$\sum \frac{\text{每个点在模拟爆发中感染的次数}}{\text{模拟爆发的总次数}}$$

某次运行所得结果如图 4.5 所示，与之对应的时间如图 4.6 所示。

```
TOP5结点影响力打分: 6.970000
TOP10结点影响力打分: 12.680000
TOP15结点影响力打分: 18.920000
TOP20结点影响力打分: 25.660000
```

图 5.5 实验结果

```
计算影响力最大的5个点用时: 43.072218s
计算这5点模拟爆发用时: 2.196754s

计算影响力最大的10个点用时: 74.938866s
计算这10点模拟爆发用时: 4.123628s

计算影响力最大的15个点用时: 129.138718s
计算这15点模拟爆发用时: 6.194436s

计算影响力最大的20个点用时: 205.958502s
计算这20点模拟爆发用时: 8.417158s
```

图 4.6 运行时间