

Influence Maximization Problem

Wang Zhiyuan 11610634

CSE

Computer Science and Technology

11610634@mail.sustc.edu.cn

1. Preliminaries

1.1. Software

In this project, I do some work to solve the Influence Maximization problem, which aid to find the initial seed set to let the final influence max in the social network.

To do this work, I do it in two part, one is ISE(Influence Size Estimator), one is IMP(Influence Maximization Problems). And the social network we use is made in two models: IC(Independent Cascade) and LT(Linear Threshold).

In the Project, I write the program by python and no extra package used. The test data is stored in the txt.

1.2. Algorithm

The algorithm I have came true is CELF and IMM. The CELF is a greedy function with pruning, and the IMM is transfer the result of the active each round to the overlay of the RR(Reverse Reachable) set.

In my test, the CELF can always get the best value in the IC model, but for the property of the LT, CELF can't always get the best result of the LT. But the largest disadvantage of the CELF is its speed, when I apply CELF to a social network with 15k nodes and 30k edges. The CELF will work more than 2 hours.

The IMM can work quickly and accurately for both IC and LT. For the graph I mentioned last paragraph, IMM can work out in 20 seconds in both IC and LT for a 50 seeds initial set with $\epsilon = 0.1$, number of processing is 8. But IMM is not perfect, this algorithm will consume so much memory when handle a graph has so many nodes with a low ϵ . For example, when I calculate a network with 425k nodes and $\epsilon = 0.1$, it need 13.2GiB memory.

2. Methodology

2.1. Representation

2.1.1. ISE. In the ISE, I do it in 3 parts:

- *BuildMap*: Read data from the file and generate the Adjacency **list**. Then reason that I choose the two dimension list but not the matrix or the dictionary is to get a balance of the memory and the speed.

- *CreateprocessingPool*: Create a processing pool to do the multiprocessing to calculate quickly.
- *DoICorLT*: Calculate the result of the network and seed given for many time.

2.1.2. CELF. In the CELF [1] [2], I do it in 3 parts:

- *BuildMap*: Read data from the file and generate the Adjacency **list**. Then reason that I choose the two dimension list but not the matrix or the dictionary is to get a balance of the memory and the speed.
- *DoISE*: Get the result of the seeds we choice in the given network and storage it in a **heap**
- *Choicenodes*: Choice that use the node in the top of the heap or do ISE again

2.1.3. IMM. In the IMM, I do it in 3 parts:

- *BuildMap*: Read data from the file and generate the Adjacency **list**. Then reason that I choose the two dimension list but not the matrix or the dictionary is to get a balance of the memory and the speed.
- *Sampling*: Calculate the influence and create the **set** of RR set.
- *NodeSelction*: Get the initial nodes set by compare the number of the RR set covered by the node.

2.2. Architecture

2.2.1. ISE.

- Read data and storage in memory.
- Get the active seeds initial.
- Active the node by the nodes actived last round.
- Do the last step until there are no new node actived in one round

2.2.2. CELF.

- Read data and storage in memory.
- Calculate the influence of the each nodes
- Choice that use the node on the top of the heap or Calculate the influence again
- Do the last step until the size of the set we get equal to the size we need

2.2.3. IMM.

- Read data and storage in memory.
- Do the sample and generate the list of the RR set
- Do node selection for the list generated last step, get the final set

2.3. Detail of Algorithm

2.3.1. ISE. In the ISE, I do 10000 times estimate in 8 processing, each do 1250 times calculations. In each calculation, the program have two parts: LT and IC, choice which parts by the parameter input. The detail of two part can look at algorithm1 and algorithm2.

algorithm 1 IC

```

1: function IC(nextNode, activeSet)
2:   activeNew  $\leftarrow$  activeSet.copy()
3:   while activeNew do
4:     activeTemp  $\leftarrow$  new set()
5:     for i in activeNew do
6:       for j in nextNode[i] do
7:         if random < j[1] then
8:           if j[0] not in activeSet then
9:             activeTemp.add(j[0])
10:            activeSet.add(j[0])
11:          end if
12:        end if
13:      end for
14:    end for
15:    activeNew = activeTemp.copy()
16:  end while return activeSet
17: end function

```

2.3.2. CELF. For the CELF, I do once ISE for each node one time and storage them in a heap depend on the influence increase of the node. Then, choice that using the node on the top of the heap or do the ISE for the nodes leftover. The detail of the CELF will be shown in the algorithm3.

2.3.3. IMM. For the IMM, it translate the ability of the activating nodes to the number of the RR set can be covered.

Defination 1 (Reverse Reachable Set). Let v be a node in network G , and g be a graph obtain by removing an each edge e in G with $1 - p(e)$ probability. The reverse reachale(RR) set for v in g is the set of nodes in g that can reach v .(That is, for each node u in RR Set, there is a directed path from u to v)

Defination 2 (Random RR set). Let \mathcal{G} be the distribution of g induced by the randomness in edge removal from G . A random RR set is an RR set generated on an instance of g randomly sampled from \mathcal{G} , for a node selected uniformly at random from g .

After the defination of the RR set and random RR set, how IMM work? By the paper Influence Maximization: Near-Optimal Time Complexity Meets Practical Efficiency [3]

algorithm 2 LT

```

1: function LT(nodes, nextNode, activeSet)
2:   threshold  $\leftarrow$  []
3:   for i  $\leftarrow$  0 to nodes do
4:     threshold.append(random())
5:   end for
6:   activeNew  $\leftarrow$  activeSet.copy()
7:   while activeNew do
8:     activeTemp  $\leftarrow$  new set()
9:     for i in activeNew do
10:      for j in nextNode do
11:        if j[0] not in activeSet then
12:          threshold[j[0]] = j[1]
13:          if threshold[j[0]] <= 0 then
14:            activeTemp.add(j[0])
15:            activeSet.add(j[0])
16:          end if
17:        end if
18:      end for
19:    end for
20:    activeNew  $\leftarrow$  activeTemp.copy()
21:  end while return activeSet
22: end function

```

algorithm 3 CELF

```

1: function CELF(nodes, size)
2:   activeSet = set()
3:   que = PriorityQueue()
4:   for i in 0 to nodes do
5:     que.add(i, (ISE(i, activeSet)))
6:   end for
7:   activeSet.add(que.get()[0])
8:   while len(activeSet) < size do
9:     B  $\leftarrow$  que.get()
10:    if ISE(B[1], activeSet) < que.get()[1] then
11:      activeSet.add(B[0])
12:    else
13:      for i in 0 to nodes do
14:        if i not in activeSet then
15:          que.add(i, (ISE(i, activeSet)))
16:        end if
17:      end for
18:    end if
19:  end while
20:  return activeSet
21: end function

```

we can know that if we add many enough RR set in the list R , and select the nodes that can cover the most number of the RR set, we must can get a good enough solution. But how many RR set we need to add in it? This is the work that IMM do more then the algorithm called TIM in the paper Influence Maximization: Near-Optimal Time Complexity Meets Practical Efficiency [4] do. First, we set some parameter that we will use in the sampling of the RR set.

$$\lambda' = \frac{(2 + \frac{2}{3}\epsilon') \cdot (\log(\frac{n}{k}) + l \cdot \log n + \log \log_2 n) \cdot n}{\epsilon'^2} \quad (1)$$

$$\lambda^* = 2n \cdot ((1 - \frac{1}{e}) \cdot \alpha + \beta)^2 \epsilon^{-2} \quad (2)$$

$$\alpha = \sqrt{l \log n + \log 2} \quad (3)$$

$$\beta = \sqrt{(1 - \frac{1}{e}) \cdot (\log(\frac{n}{k}) + l \log n + \log 2)} \quad (4)$$

First of all, let's look at the algorithm the NodeSelection, in this part, we will get a given size set from the nodes in the G that let the nodes in the set can cover as more as possible. This part not only use in the last step to get the target set, but also use in the sampling part to determine the size of the RR set list we need. In the NodeSelection part,

algorithm 4 NodeSelection

```

1: function NODESELECTION( $\mathcal{R}, k$ )
2:    $S_k^* \leftarrow \text{set}()$ 
3:   for  $i \leftarrow 1$  to  $k$  do
4:      $v \leftarrow \text{Maximizes } \text{num}(\text{RR Set Covered})$ 
5:      $S_k^*.add(v)$ 
6:   end for
7: end function
8: return  $S_k^*$ 

```

how to let the speed of line 4 be highest is the key of the optimization of the algorithm. When I implment it, I use a **map** to storage many **list** as the **value**, each **list** storage the nodes that can be reached by the **key** node of this **list**. And at the same time, I count that the number of the nodes that can be reached for each **key** node Then, I will do the work that choice the node to add to the target set. I get the node that have the max number of the nodes can be reached, then update the map and the list that used to statistic last step. I will repeat this step until the target set is full.

Then, I will do the Sampling part work. The detail of the sampling is in the algorithm5. In the sampling, it calculate the number of the RR set needed Dynamicly and can save much time.

3. Empirical Verification

3.1. Design

When I do the verification. I set the *hyperparameters* $\epsilon = 0.1$, $l = 1$.

algorithm 5 Sampling

```

1: function SAMPLING( $G, k, \epsilon, l$ )
2:    $\mathcal{R} \leftarrow []$ 
3:    $LB \leftarrow 1$ 
4:    $\epsilon' \leftarrow \sqrt{2} \cdot \epsilon$ 
5:   for  $i \leftarrow 1$  to  $\log_2 n - 1$  do
6:      $x \leftarrow n/2^i$ 
7:      $\Theta_i \leftarrow \lambda'/x$ 
8:     while  $|\mathcal{R}| \leq \Theta_i$  do
9:        $v \leftarrow \text{random}(G.\text{node})$ 
10:       $\mathcal{R}.append(v)$ 
11:    end while
12:     $S_i \leftarrow \text{NodeSelection}(\mathcal{R}, k)$ 
13:    if  $n \cdot \mathcal{F}_{\mathcal{R}}(S_i) \geq (1 + \epsilon' \cdot x)$  then
14:       $LB \leftarrow n \cdot \mathcal{F}_{\mathcal{R}}(S_i) / (1 + \epsilon')$ 
15:    break
16:  end if
17: end for
18:  $\Theta \leftarrow \lambda^* / LB$ 
19: while  $|\mathcal{R}| \leq \Theta$  do
20:    $v \leftarrow \text{random}(G.\text{node})$ 
21:    $\mathcal{R}.append(v)$ 
22: end while
23: return  $\mathcal{R}$ 
24: end function

```

3.1.1. Network. : NetHEPT, Amazon and DBLP

3.1.2. Hyperparameter. : $\epsilon = 0.1$ or 0.2 , $l = 1$. The ϵ less, the speed lower, but the accuracy higher.

3.2. Performance

I test the program on my computer with AMD Ryzen 17003.6GHz@16 and 16GB memory And the test for the NetHEPT use 16 processing and the test of the Amazon and DBLP will use 1 processing becauseing The memory will not enough when calculate a very large network

- N: Num of nodes; E: Num of edges; M: Type of model
- S: Size of the seeds; P: Num of processing
- T: Time used; R: The result

TABLE 1. RESULT OF TEST CASE

Network	N	E	M	S	ϵ	P	T	R
NetHEPT	15233	32235	LT	50	0.1	16	13.83	1701.97
NetHEPT	15233	32235	IC	50	0.1	16	15.71	1295.52
NetHEPT	15233	32235	LT	50	0.2	16	4.33	1698.73
NetHEPT	15233	32235	IC	50	0.2	16	4.56	1290.56
Amazon	548526	13474	LT	50	0.2	1	189.55	852.99
Amazon	548526	13474	IC	50	0.2	1	347.62	636.59
Amazon	548526	13474	LT	50	0.5	1	70.15	847.83
Amazon	548526	13474	IC	50	0.5	1	69.76	628.03
DBLP	425746	15794	LT	50	0.1	1	376.17	1002.82
DBLP	425746	15794	LT	50	0.2	1	168.81	1001.21
DBLP	425746	15794	IC	50	0.2	1	181.13	658.26

3.3. Result

From the table above, we can find that IMM have good effect to the IMP problem, but it work slowly in very large network with small ϵ .

3.4. Analysis

For the test, I find that the IMM can work fast, but the performance is limited by the speed of node selction because it can not run in multiprocessing, so the next optimization step is add nodeselection into the multiprocessing. Another limitation of IMM is the memory, when a large number of nodes calculate with a low ϵ , the number of RR set needed will be so large, and these RR set is required in the last node selection. So it will take so much memory, and the multiprocessing also can not be used.

References

- [1] J. Lv, J. Guo, Z. Yang, W. Zhang, and A. Jocshi, "Improved algorithms of celf and celf++ for influence maximization," *Journal of Engineering Science & Technology Review*, vol. 7, no. 3, 2014.
- [2] A. Goyal, W. Lu, and L. V. Lakshmanan, "Celf++: optimizing the greedy algorithm for influence maximization in social networks," in *Proceedings of the 20th international conference companion on World wide web*. ACM, 2011, pp. 47–48.
- [3] Y. Tang, Y. Shi, and X. Xiao, "Influence maximization in near-linear time: A martingale approach," in *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data*. ACM, 2015, pp. 1539–1554.
- [4] Y. Tang, X. Xiao, and Y. Shi, "Influence maximization: Near-optimal time complexity meets practical efficiency," in *Proceedings of the 2014 ACM SIGMOD international conference on Management of data*. ACM, 2014, pp. 75–86.