

# Influence Maximization Problem

YingZhou 11610701  
Computer Science and Engineering  
SUSTech  
11610701@mail.sustc.edu.cn

## 1. Preliminaries

### 1.1. Problem Description

The *influence maximization problem* (IMP) normally asks for  $k$  finding a small subset of  $k$  nodes (referred to as seed set) in a social network  $G$  that could maximize the spread of influence which is the expected number of nodes that are influenced by the nodes in the seed set in a cascade manner (ISE). The project includes two parts which are influence speed estimate (ISE) and influence maximization problem (IMP).

### 1.2. Problem Application

The *influence maximization problem* (IMP) has been extensively studied recently due to its potential commercial value. A well-known application is the case that a company wants to spread the adoption of a new product from some initially selected adopters through the social links between users.

## 2. Methodology

### 2.1. Notation

TABLE 1: Notation

Notation	Description
$G = (V, E)$	a social network $G$ with a node set $V$ and an edge set $E$
$n, m$	the numbers of nodes and edges in $G$ , respectively
$k$	the size of the seed set for influence maximization
$R$	the set of RR sets generated by IMM's sampling phase
$\theta$	the number of RR sets in $R$
$F_R(s)$	the fraction of RR sets in $R$ that are covered by a node set $S$

### 2.2. Data Structure

The main data structures I used contain dictionary, queue, set and heap et. Dictionary is used to store the edges and the weight of each edge since using adjacency matrix is memory wasting if all graphs given are sparse. Queue is applied in multi-processes part for workers to get and

put the jobs. Set replacing list records active vertices to save initializing time because set in Python is implemented in hash table. Finally, lazy update optimization with heap implemented in node selection part.

### 2.3. Model Design

Both ISE and IMP are accelerated with 8 processes. ISE is strictly implemented according to the pseudo code of tutorial lecture. The novel point about IMM is that, in a social network, a influential person started by many people may not follow same number of people. In that case, the sub-graph of converse graph is much smaller than that of original graph. It is simple and straightforward in implementation process as the description of the paper [1] is clear and specific.

### 2.4. Detail of Algorithm

**2.4.1. Influence Spread Estimate.** Given the seed set  $S$ , the diffusion processes unfolds in discrete rounds as follow: At round 0, all nodes in  $S$  are active and the others are inactive. In the subsequent rounds, the newly activated nodes will try to activate their neighbors. The process will stop when no more nodes get activated. There are two fundamental diffusion models are specified in this project including *Independent Cascade* (IC) and *Linear Threshold* (LT).

**Independent Cascade** When a node  $u$  gets activated, initially or by another node, it has a single chance to activate each inactive neighbor  $v$  with the probability proportional to the edge weight  $w(u, v)$ .

**Linear Threshold** At the beginning, each node  $v$  selects a random threshold  $\theta_v$  uniformly at random in range  $[0, 1]$ . Since the second round, an inactive node  $v$  becomes activated if  $\sum \text{activated neighbors } u W(u, v) \geq \theta_v$ . The weight of the edge  $(u, v)$  equals  $\frac{1}{d_{in}(v)}$ .

---

**Algorithm 1** Independent Cascade

---

```

1: function IC
2:   Initialize ActivitySet as {SeedSet}
3:   count  $\leftarrow$  ActivitySet.length
4:   while ActivitySet is not empty do
5:     Set newActivitySet as  $\emptyset$ 
6:     for each seed in ActivitySet do
7:       for each inactive neighbor in seed do
8:         Seed tries to activate neighbors with
probability
9:         if activated then
10:          Add neighbor into activate to newAc-
tivitySet
11:        end if
12:      end for
13:    end for
14:    count  $\leftarrow$  count + newActivitySet.length
15:    ActivitySet  $\leftarrow$  newActivitySet
16:  end while
17: end function
18: return count

```

---



---

**Algorithm 2** Linear Threshold

---

```

1: function LT
2:   Initialize ActivitySet as {SeedSet}
3:   Randomly sample thresholds
4:   count  $\leftarrow$  ActivitySet.length
5:   while ActivitySet is not empty do
6:     Set newActivitySet as  $\emptyset$ 
7:     for each seed in ActivitySet do
8:       for each inactive neighbor in seed do
9:         Calculate the weights of activated neigh-
bors as w_total
10:        if w_total  $\geq$  neighbor.thresh then
11:          Update the state as Active
12:          newActivitySet
13:        end if
14:      Add neighbor into activate to newAc-
tivitySet
15:    end for
16:  end for
17: end while
18: end function

```

---

**2.4.2. Influence Maximization Problem.** All the details and proofs about the algorithm I used in IMP are referred Y.Tang’s work [1]. I just implemented it and decided which data structures to use to make code run faster.

---

**Algorithm 3** IMM

---

```

1: function IMM
2:    $l = l \cdot (1 + \log(2) / \log(n))$ 
3:    $R = \text{Sampling}(G, k, \epsilon, l)$ 
4:    $S_k = \text{NodeSelection}(R, k)$ 
5: end function

```

---



---

**Algorithm 4** Sampling

---

```

1: function SAMPLING( $G, k, \epsilon, l$ )
2:   Initialize a set  $R = \emptyset$  and an integer  $LB = 1$ 
3:   Let  $\epsilon l = \sqrt{2} \cdot \epsilon$ 
4:   for  $i = 1$  to  $\log_2 n$  do
5:      $x = n/2^i$ 
6:      $\theta_i = \lambda l / x$ , where  $\lambda l$  is as defined in equation 1
7:     while  $\|R\| \leq \theta_i$  do
8:       Select a node  $v$  from  $G$  uniformly at random
9:       Generate an RR set for  $v$ , and insert it into
 $R$ 
10:    end while
11:    Let  $S_i = \text{NodeSelection}(R)$ 
12:    if  $n \cdot F + R(S_i) \geq (1 + \epsilon l) \cdot x$  then
13:       $LB = n \cdot F_R(S_i) / (1 + \epsilon l)$ 
14:    break
15:  end if
16: end for
17: Let  $\theta = \lambda^* / LB$ , where  $\lambda^*$  is as defined as in
Equation 2
18: end function

```

---



---

**Algorithm 5** NodeSelection

---

```

1: function NODESELECTION( $R, k$ )
2:   Initialize a node set  $S_k = \emptyset$ 
3:   for  $i = 1$  to  $k$  do
4:     Identify the vertex  $v$  that maximizes  $F_R(S_k \cup$ 
 $v) - F_R$ 
5:     Insert  $v$  into  $S_k$ 
6:   end for
7: end function

```

---

### 3. Empirical Verification

#### 3.1. Dataset

TABLE 2: Datasets

Name	$n$	$m$	Type
NetHEPT	15.2K	32.2K	undirected
Epinion	39.0M	84.0K	directed
Twitter	81.3K	1.8M	undirected

#### 3.2. Hyperparameters

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

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

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

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

Most parameters are proposed in Y.Tang's work [1]. The default setting of Algorithm 3 is  $\epsilon = 0.1$  and  $l = 1$ . The maximization algorithm the paper presented provides a  $(1 - 1/2 - \epsilon)$ -approximate solution with at least  $1 - 1/n^l$  probability. The  $\epsilon$  and  $l$  setting is a trade-off between solution quality and time cost, and is set as the same as most people.

### 3.3. Performance Measure

The performances of accuracy and efficiency are considered. I measure the accuracy of results by comparing time and results with best performance in OJ and the efficiency presented by running in larger instances. The code is written in Python, compiled using Python 3.6.5 :: Anaconda. Only NumPy package is extra imported. Ubuntu 18.04.1, 16 processors, each processor with 2 threads, 32G RAM, Intel(R) Xeon(R) CPU E5-2620 v4 @ 2.10GHz.

### 3.4. Experimental Result

TABLE 3: The performance in NetHEPT instance ( $\epsilon = 0.2$ )

k	diffusion model	time_avg(s)	ISE	time_OJ	OPT_OJ
5	IC	1.32	324.16	10.04	324.16
	LT	1.04	393.96	2.557	392.98
50	IC	1.37	1291.89	53.98	1298.10
	LT	1.21	1699.58	8.398	1702.00

TABLE 4: The performance in NetHEPT instance ( $\epsilon = 0.1$ )

k	diffusion model	time_avg(s)	ISE	time_OJ	OPT_OJ
5	IC	3.34	323.31	10.04	324.16
	LT	2.78	393.72	2.557	392.98
50	IC	3.99	1295.20	53.98	1298.10
	LT	2.91	1701.19	8.398	1702.00

TABLE 5: The performance in NetHEPT instance ( $\epsilon = 0.05$ )

k	diffusion model	time_avg(s)	ISE	time_OJ	OPT_OJ
5	IC	10.04	323.72	10.04	324.16
	LT	7.87	393.72	2.557	392.98
50	IC	12.34	1297.22	53.98	1298.10
	LT	8.52	1702.00	8.398	1702.00

TABLE 6: The performance in larger instances ( $\epsilon = 0.1$ )

Instance	k	diffusion model	time(s)
Epinion	5	IC	6.25
		LT	4.17
	50	IC	8.75
		LT	6.72
Twitter	5	IC	55.12
		LT	15.41
	50	IC	70.51
		LT	16.70

### 3.5. Conclusion

According to the experiment results, it proves that my code can get the same quality solutions in less time. To conclude, the algorithm proposed in [1] is creative and efficient. Not only reversing graph to solve *influence maximization problem* is proposed, but also it calculates the lower bounds of sampling size saving computation cost. The obvious disadvantage is that it still needs massive sampling which is memory consuming. This project greatly helps me to understand the extremely important role of data structure and code construction when I completed some tricky idea in my code. If there is a chance, I would like to contribute some works in *influence maximization* problem.

### References

- [1] Y. Tang, Y. Shi, and X. Xiao, Influence Maximization in Near-Linear Time, in Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data - SIGMOD 15, 2015.