

# Influence Maximization Problem Project

Wentao Ning 11610914  
Computer Science and technology  
Southern University of Science and Technology  
11610914@mail.sustc.edu.cn

## 1. Preliminaries

**Description:** Influence Maximization Problem (IMP) is the problem of finding a small subset of nodes (referred to as seed set) in a social network that could maximize the spread of influence. The influence spread is the expected number of nodes that are influenced by the nodes in the seed set in a cascade manner. The Linear Threshold and Independent Cascade Models are two of the most basic and widely-studied diffusion models[1]

**Application:** The IMP is often used to calculate the influence of a social network and so on. A social network, the graph of relationships and interactions within a group of individuals, plays a fundamental role as a medium for the spread of information, ideas, and influence among its members[1]

### 1.1. Software

I wrote this project in Python 3.6 using editor Sublime Text 2 and IDE PyCharm.

### 1.2. Algorithm

At first I wrote a preprocess function to preprocess the file and get the node number and all edges. Then I used CELF algorithm which is a kind of greedy search algorithm based on submodularity to get the final result.

## 2. Methodology

### 2.1. Representation

There are 3 global variables in both ISE.py and IMP.py.

- **graph\_list:** A dict which records the adjacent list of each node
- **graph\_parent\_list:** A dict which records the parent nodes of each node
- **neighbor\_prob:** A dict which records the weight of each node

### 2.2. Architecture

There are several methods in my ISE.py file.

- **preprocess:** Preprocess the given file and construct the graph.
- **calculate\_influence:** To calculate the influence based on different diffusion model.
- **one\_IC\_sample:** To calculate the influence based on IC diffusion model.
- **one\_LT\_sample:** To calculate the influence based on LT diffusion model.

There are some other functions in my IMP.py file.

- **CELF:** Using a kind of greedy algorithm called CELF to calculate the result set.
- **print\_results:** To print result as the required format.

### 2.3. Detail of Algorithm

There are 3 main functions in my ISE.py file.

- **one\_IC\_sample:** To calculate the influence based on IC diffusion model.

---

**Algorithm 1** one\_IC\_sample

---

```
1: ActivitySet  $\leftarrow$  SeedSet
2: count  $\leftarrow$  len(ActivitySet)
3: while ActivitySet is not empty do
4:   newActivitySet  $\leftarrow$   $\emptyset$ 
5:   for each seed in ActivitySet do
6:     for each inactive neighbor in seed do
7:       seed tries to activate neighbors with probability
8:       if activated then
9:         newActivitySet.add(neighbor)
10:      end if
11:    end for
12:  end for
13:  count  $\leftarrow$  count + len(newActivitySet)
14:  ActivitySet  $\leftarrow$  newActivitySet
15: end while
16: return count
```

---

- **one\_LT\_sample:** To calculate the influence based on LT diffusion model.

---

**Algorithm 2** one\_LT\_sample

---

```

1: ActivitySet  $\leftarrow$  SeedSet
2: Randomly set threshold for each node
3: count  $\leftarrow$  len(ActivitySet)
4: while ActivitySet is not empty do
5:   newActivitySet  $\leftarrow$   $\emptyset$ 
6:   for each seed in ActivitySet do
7:     for each inactive neighbor in seed do
8:       calculate the weight of activated neighbors: w_total
9:       if w_total  $\geq$  neighbor.threshold then
10:        update the state as Active
11:        newActivitySet.add(neighbor)
12:      end if
13:    end for
14:  end for
15:  count  $\leftarrow$  count + len(newActivitySet)
16:  ActivitySet  $\leftarrow$  newActivitySet
17: end while
18: return count

```

---

- **calculate\_influence:** To calculate the influence based on 10000 times iterations on LT or IC diffusion model.

---

**Algorithm 3** calculate\_influence(seed, model)

---

```

1: sum  $\leftarrow$  0
2: times  $\leftarrow$  10000
3: if model.upper() = "IC" then
4:   for i = 1 to times do
5:     sum  $\leftarrow$  sum + one_IC_sample(seeds)
6:   end for
7:   return sum / times
8: else if model.upper() = "LT" then
9:   for i = 1 to times do
10:    sum  $\leftarrow$  sum + one_LT_sample(seeds)
11:   end for
12:   return sum / times
13: end if

```

---

There is one main function in my IMP.py file.

**CELLF:** Using a kind of greedy algorithm called CELF to calculate the result set.

---

**Algorithm 4** CELF

---

```

1: influence_dict  $\leftarrow$  empty dict
2: for i = 1 to len(neighbor_prob) do
3:   influence_dict[i]  $\leftarrow$  calculate_influence(i, model)
4: end for
5: seeds  $\leftarrow$  empty set
6: tmp_max_node  $\leftarrow$  the key in influence_dict having the max value
7: add tmp_max_node2 to seeds
8: pop tmp_max_node from influence_dict
9: while len(seeds) < seed_num do
10:  tmp_max_node  $\leftarrow$  the key in influence_dict having the max value
11:  influence_dict[tmp_max_node]  $\leftarrow$  the influence increment after adding tmp_max_node
12:  tmp_max_node2  $\leftarrow$  the key in influence_dict having the max value
13:  if tmp_max_node2 = tmp_max_node then
14:    add tmp_max_node2 to seeds
15:    pop tmp_max_node2 from influence_dict
16:  else
17:    continue
18:  end if
19: end while
20: return seeds

```

---

### 3. Empirical Verification

#### 3.1. Design

In the Empirical Verification part, I use the network.txt which is given by teacher as the input graph and the seeds.txt and seeds2.txt as input seed set to verify the ISE model. Then I use network.txt as input to verify the IMP model.

#### 3.2. Data and data structure

- **Data Set:** NetHEPT.txt, network.txt, seeds.txt and seeds2.txt
- **Data Structure:** list, set, dict and numpy 2D matrix.

#### 3.3. Performance & Hyperparameters

**Performance Measurement:** I run my algorithm on the given files, network.txt and NetHEPT.txt. I test my ISE.py file to calculate the correctness using seeds.txt and seeds2.txt file. I test my IMP.py file to get 5 seeds and 10 seeds respectively.

**Hyperparameters setting:** At first, I set the ISE Algorithm's iteration times to 10000 in my IMP.py file, When I run the IMP.py file on the NetHEPT.txt, I found it is too slow and can not get a good result on the given time. So I change the iteration times to 1000, I found the result is not change but the run time decreased dramatically. Then I

changed it to 500, the result is also unchanged but the run time decreased half. Then I continue decreasing the iteration times, but this time when I decrease the iteration times, the performance of this algorithm will be worse. So I set the iteration times to 500 finally.

### 3.4. Result

In the ISE test part, I let the one\_IC\_sample or one\_LT\_sample iterate 10000 times and recorded the result and used time. The result is below.

	IC	LT
<b>seeds.txt</b>	5.019 / 0.23s	5.033 / 0.05s
<b>seeds2.txt</b>	30.47 / 0.25s	37.07 / 0.64s

TABLE 1. INFLUENCE RESULT AND USED TIME OF EACH DIFFUSION MODEL AND SEED SET

In the IMP test part, I set the seed set as 5 and 10 respectively, the result and the used time is below.

	IC	LT
<b>5</b>	30.674 / 1.959s	30.719 / 1.831s
<b>10</b>	42.75 / 5.587s	42.862 / 6.616s

TABLE 2. INFLUENCE RESULT AND USED TIME OF EACH DIFFUSION MODEL AND SEED SIZE

### 3.5. Conclusion

This project is to solving a problem called Influence Maximum Problem. I used a kind of greedy algorithm to solve this problem. Admittedly, the general greedy algorithm is not effective, so I used an algorithm called CELF which consists of greedy algorithm plus pruning based on submodularity. The advantage of my algorithm is that this algorithm can have a good performance on the graph which is not too large, and it is very fast to get the results. In my opinion, the disadvantage of my algorithm is that my algorithm can not get a good performance on the large graph, it will take several minutes or hours to run on the large graph, which is what I need to improve.

### Acknowledgments

### References

- [1] D. Kempe, J. Kleinberg and E. Tardos, "Maximizing the Spread of Influence through a Social Network", Cs.cornell.edu, 2003. [Online]. Available: <https://www.cs.cornell.edu/home/kleinber/kdd03-inf.pdf>. [Accessed: 01- Dec- 2018].