

IMP Report

YILIN ZHANG 11510280

School of Computer Science and Engineering
Southern University of Science and Technology
11510280@mail.sustc.edu.cn

December 15, 2018

1. PRELIMINARIES

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.

A social network is modeled as a directed graph $G = (V, E)$ with nodes in V modeling the individual in the network and each edge $(u, v) \in E$ is associated with a weight $w(u, v) \in [0, 1]$ which indicates the probability that u influences v . Let $S \subseteq V$ to be the subset of nodes selected to initiate the influence diffusion, which is called the seed set. The stochastic diffusion models specify the random process of influence cascade from S , of which the output is a random set of nodes influenced by S . The expected number of influenced nodes $\sigma(S)$ is the influence spread of S . And the goal of IMP is to find a seed set S that maximizes $\sigma(S)$, subject to $|S| = k$.

In this project, we need to implement an Influence Spread Estimator (ISE) and an IMP solver named ISE.py and IMP.py, respectively. The first is to implement an estimation algorithm for the influence spread and the second is to design and implement a search algorithm for IMPs. The IMP is NP-hard and the influence spread computation is #P-hard under the definitions shown in the above.

1.1. Software

This project is written in Python. The libraries being used includes *argparse*, *sys*, *Numpy* and *Random*. The local test environment is as follows:

- Python version: 3.6.6

- Operation System: macOS
- Processor: 3.5 GHz Intel Core i7
- Memory: 16 GB

2. METHODOLOGY

Given the seed set, the diffusion process unfolds in discrete rounds:

- 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.
- Once a node becomes active, it will remain active till the end.
- The process stops when no more nodes get activated.

Two fundamental diffusion models Linear Threshold (LT) and Independent Cascade (IC) are specified in this project. I estimate the influence spread by these two models, and search seeds by Cost-effective Lazy-forward (CELF) method.

2.1. Representation

Both ISE and IMP problem need to parse arguments from command line and read *.txt files including the information of graph or seeds. I used a dictionary to store general informations such as **NETWORK_FILE_PATH**, **DIFFUSION_MODEL**, **TERMINATION** etc. As for the graph, I used a 2-dimensional matrix denotes the influence weights among them.

The format of solvers call are as follows (the test environment is Unix-like):

Task 1:

python ISE.py -i <social network> -s <seed set> -m <diffusion model> -t <time budget>

Output: the value of the estimated influence spread.

Task 2:

python IMP.py -i <social network> -k <predefined size of the seed set> -m <diffusion model>

Output: the seed set found by the algorithm.

- **<social network>**: absolute path of the social network file.
- **<seed set>**: absolute path of the seed set file.
- **<predefined size of the seed set>**: a positive integer.
- **<diffusion model>**: IC or LT.
- **<time budget>**: a positive number which indicates how many seconds the algorithm can spend on this instance.

Some constants are as follows:

- **INFO**: a dict stores informations from args.
- **graph**: two-dimensional matrix to present the directed graph.
- **precision**: determine estimate times.

2.2. Architecture

The executable files are *ISE.py* and *IMP.py*, which contain parameter parsing, file reading, and result printing. There are four modules *IC.py*, *LT.py*, *CELF.py* and *utils.py* in the *submodule* folder, offer functions IC, LT, CELF and other tools, respectively.

- **IC**:
 - **IC**: This program is to estimate the influence, and output the value of the estimated influence spread. 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 $w(u, v)$. Afterwards, the activated nodes remain its active state but they have no contribution in later activations. After there is no new activated node, return the count of activated nodes as influence.
- **LT**:

- **LT**: At the beginning, each node v selects a random threshold θ_v uniformly at random in range $[0,1]$. If round $t1$, an inactive node v becomes activated if for all activated neighbors u , $\sum w(u, v) \geq \theta_v$. After there is no new activated node, return the count of activated nodes as influence.

- **CELF**:

- **influence**: Calculate the marginal influence and return it.
- **CELF**: Based on greedy algorithm and exploiting submodularity, significantly reduce # influence spread evaluations. Return the set of seeds.

- **utils**:

- **count**: Calculate k times seed set influence and return the average value.
- **model_adapter**: Transfer diffusion model from string to int to simplified judgment.

Algorithm 1 IC

Input: *graph, SeedSet*

Output: count of activated nodes (influence)

```

1: initial ActivitySet = SeedSet
2: count = ActivitySet.length
3: while !ActivitySet.IsEmpty() do
4:   initial newActivitySet =  $\emptyset$ 
5:   for all seed  $\in$  ActivitySet do
6:     for all inactive neighbor  $\in$  seed do
7:       try to activate by  $w(\text{seed}, \text{neighbor})$ 
8:       if activate succeeded then
9:         update neighbor's state
10:        newActivitySet.add(neighbor)
11:      end if
12:    end for
13:  end for
14:  count = count + newActivitySet.length
15:  ActivitySet = newActivitySet
16: end while
17: return count

```

2.3. Algorithm

ISE program simply use IC or LT model to estimate the influence. It prefer an average value

Algorithm 2 LT

Output: count of activated nodes (influence)

```

1: initial ActivitySet = SeedSet
2: for all node do
3:   initial threshold  $\in [0, 1]$ 
4:   if threshold = 0 then
5:     add node to ActivitySet
6:   end if
7: end for
8: count = ActivitySet.length
9: while !ActivitySet.IsEmpty() do
10:  initial newActivitySet =  $\emptyset$ 
11:  for all seed  $\in$  ActivitySet do
12:    for all inactive neighbor  $\in$  seed do
13:      calculate the sum wtotal of all ac-
      tive neighbors' weight of neighbor
14:      if wtotal  $\geq$  neighbor.threshold then
15:        update neighbor's state
16:        newActivitySet.add(neighbor)
17:      end if
18:    end for
19:  end for
20:  count = count + newActivitySet.length
21:  ActivitySet = newActivitySet
22: end while
23: return count

```

by repeatedly estimated when time is limited, so we define a number by **formula 1** named *precision*, which is inversely proportional to the graph size. The result need to plus 1 to avoid value 0, and have an upper bound *MAX_{precision}* to guarantee execution efficiency.

IMP program also used the *precision*, but it has different meaning with it in ISE program. IMP program using CELF algorithm to get seeds. CELF function just running once in this process, but for estimating emarginal influences (Δinf) **formula 2**, average assessment is essential to guarantee the reliability of each single seed choice.

Here are details of these algorithms.

$$precision = MIN(\frac{scale_factor}{node_number} + 1, MAX_{precision}) \quad (1)$$

Algorithm 3 influence

Input: *graph, current_seed, model, precision*

Output: current influence of *current_seed*

```

1: initial infs = 0
2: if model is IC then
3:   repeat
4:     infs = infs + IC(graph, current_seed)
5:   until precision times
6: else
7:   repeat
8:     infs = infs + LT(graph, current_seed)
9:   until precision times
10: end if
11: infs = infs / precision
12: return infs

```

$$node_{\Delta inf} = influence_{SeedSet \cup node} - influence_{SeedSet} \quad (2)$$

3. EMPIRICAL VERIFICATION

3.1. Experiments

A basic test (a graph with 62 nodes and 159 edges) has been given to test the usability of ISE and IMP code, and both of them running well. Also, my code can handle further tests on online-judge platform in limited time and get good solutions except LT model of the bonus dataset named NetHEPT-50 which has more than 15,000 nodes. But considering I just using CELF rather than more advanced algorithm, this result indicates that the code has a good performance.

3.2. Data and data structure

Data being used in this project is a directed graph and the influence number between its nodes. Data structure used in this project includes dictionary, array, tuple, list, and set.

3.3. Performance

CELF's theoretic running time, computed by time complexity of functions is about $O(krn'\sigma)$, where *k* is size of seeds, σ is the influence spread

Algorithm 4 CELF

Input: $graph, size, model, precision$ **Output:** a k size seed set

```
1: initial  $SeedSet = \emptyset$ 
2: initial  $queue = \emptyset$ 
3:  $current\_influence = 0$ 
4: add all nodes with nonzero out-degree to
   queue
5: for all  $node_i \in queue$  do
6:   initial  $\Delta inf_i = \infty$ 
7: end for
8: while  $SeedSet < size$  do
9:   initial  $best = 0$ 
10:   $new\_queue =$  reverse sorted  $queue$  by
     marginal influence
11:  while exit  $\Delta inf_{node} \in new\_queue > \Delta inf_{best}$ 
     do
12:    update  $\Delta inf_{node}$ 
13:    if  $n$  then  $\Delta inf_{node} > \Delta inf_{best}$ 
14:       $best = node$ 
15:    end if
16:  end while
17:  add  $best$  to  $SeedSet$ 
18:  update  $current\_influence$ 
19: end while
20: return  $SeedSet$ 
```

REFERENCES

- [1] Yuchen Li, Ju Fan, Yanhao Wang, and Kian-Lee Tan, *Influence Maximization on Social Graphs: A Survey*, IEEE Transactions on Knowledge and Data Engineering, pp. 1852-1872, 2018.

of the seed set. n' is the number of simulations at each iteration.

3.4. Analysis

IC and LT are two of classical models to analyze IMP problem, and there are some other models are not considered in this project such as the triggering model (TR) and the Continuous-Time (CT) IC model[1] *etc.* Besides, CELF is a simulation method mentioned in 2007. There are some better approaches now considering IMP problem by improving *theoretical efficiency* of the simulation based methods while preserving the approximation guarantee, these new approaches are able to handle bigger and more complex network graph than CELF.