

Project 4: Support Vector Machine

Yuejian Mo 11510511

Department of Biology

Southern University of Science and Technology

Email: 11510511@mail.sustc.edu.cn

1. Preliminaries

Support Vector Machine(SVM) is... Our data ...

Support vector machine, are supervised learning models with associated learning algorithms that analyze data used for classification and regression analysis. Given a set of training examples, each marked as belonging to one or the other of two categories, an SVM training algorithm builds a model that assigns new examples to one category or the other, making it a non-probabilistic binary linear classifier. [?] If the dataset are points in space, we can find out an gap separate dataset into two part by using SVM.

This project is an implementations of Greedy algorithm to find out a seed set to maximize the spread of influence through a social network, based on LT(Linear Threshold) and IC(Independent Cascade) models. With the development of Internet and smart phone, more and more people are connected by social network. It is interesting to study the influence spread between different members. Explicitly, maximizing the spread in network with less afford takes the eyes of adversary provider and government managers. In 2003, Kempe, Kleinberg and Tardos proposal two models to represent and optimize this kind of problem. [?]

This problem, Influence Maximizing Problem(IMP) is a NP-hard. Inspired by above report, I success produce seed set by natural greedy.

In this problem, I use SVM with random gradient down to classify training data successful.

1.1. Software

This project is written by Python 3.7 with editor Vim. Numpy, os, time, random, sys and argparse library are used.

1.2. Algorithm

Model question with SVM. To optimal time cost and current of solution, random gradient down is used.

2. Methodology

Training data is
SVM is modeled as ...
random

Firstly, a social network is modeled as a directed graph $G = (V, E)$ and each edge $(u, v) \in E$ is associated with a weight $w(u, v) = \frac{1}{d_{in}^v}$ which indicates the probability that u influences v . $S \in V$ is the subset of nodes selected to initiate the influence diffusion, which is called seeds set.

Then, I try to evaluate the spread influence with two different diffusion models. In *Linear Threshold* model, a node v is influenced by each neighbor u according to a weight $w_{v,u}$ such that $\sum_u w_{v,u} \leq 1$. The dynamics is following. Each node v are accessed a threshold $\theta \in [0, 1]$ randomly obeyed uniform distribution at first step; this represents the weighted fraction of v ' neighbors that must become active in order for v become active. In step t , each inactive node v will be active by its active neighbors w at step $t - 1$ as following

$$\sum_w w_{w,v} \geq \theta_v$$

. Once no new active node is generated, maximal spread influence will obtain with given seeds set in specified graph. I notices final influence as the total number of active node σ .

In *Independent Cascade* model, I start with an initial set of active nodes S , and the process unfolds in discrete steps according to the following randomized rule. When node v first becomes active in step t , it is given a single chance to activate each currently inactive neighbor w ; it succeeds with a probability $w_{v,w}$, independently of the history thus far. Every step. In practical, for each inactive neighbor w of new active node v will be active, if $w_{v,w}$ is larger than a new generated random number between 0 and 1 obeyed uniform distribution. (If w has multiple newly activated neighbors, their attempts are sequenced in an arbitrary order.) If v succeeds, then it cannot make any further attempts to activate w in subsequent rounds. Again, the process runs until no more activation are possible.

Both of Linear Threshold and independent Cascade model are stimulated by repeating many time to average influence for each seed.

Finally, the seed required for IMP is generated by a natural greedy algorithm. Each time, a node with highest influence will added to a seeds set. This empty seed set will fill with k node, where k is given from outside.

2.1. Representation

Some main data are maintain during process: **time_budget**, **node_num**, **graph_cp**, **graph_pc**, **incoming**. Others data would be specified inside functions.

- **time_budget**: How many second (in Wall clock) spend on instance, which provided from OJ system.
- **node_num**: The total node number of given graph.
- **graph_cp**: The adjacent list of network graph using dictionary, where a node point to all precursors, which is generated from input file.
- **graph_pc**: The adjacent list of network graph using dictionary, where a node point to all descendants, which is generated from input file.
- **incoming**: A dictionary stores the incoming weight $w_{v,u}$, indexed by node u .

2.2. Architecture

Here list all functions of **ISE.py** in given code:

- **genSeeds**: Read seeds and store in a list from input file.
- **genGraph**: Generate *graph_cp* and *graph_pc* from input file.
- **LT**: Calculate influence of linear threshold model of given seeds set.
- **IC**: Calculate influence of independent cascade model of given seeds set.
- **calSpread**: Calculate influence by repeating many time of *LT* or *IC* function.
- **s_format**: Standard output function.
- **__main__**: Main control function.

The **IMP.py** is similar to **ISE.py**, but one function is added:

- **hill_greedy**: Find out optimal seeds set using Hill-Greedy algorithm.

The **ISE.py** and **IMP.py** were executed in OJ system and local system.

2.3. Detail of Algorithm

Here describes some vital functions.

- **genGraph**: Generate global cost and demand graph from input file.

Algorithm 1 genGraph

Input: *input_file_name*

Output: *graph_cp, graph_pc, incoming*

- 1: open *input_file_name* as *file* {open file and read line by line}
 - 2: *node_num, edge_num* \leftarrow first line of *file*
 - 3: {Read each line information until arrive edge information}
 - 4: **for** each line in *file* **do**
 - 5: *precursor, descendant, weight* \leftarrow split each line
 - 6: *graph_cp*[*precursor*] = *descendant*
 - 7: *graph_pc*[*descendant*] = *precursor*
 - 8: *incoming*[*descendant*] = *weight*
 - 9: **end for**
-

- **LT**: Calculate the influence of seeds set of Linear Threshold model.

Algorithm 2 LT

Input: *graph_pc, S, incoming*

Output: length of *isActive*

- 1: copy *S* to new list *isActive*
 - 2: create empty list *threshold*
 - 3: **for** each node in graph **do**
 - 4: Add Random Number α between [0,1] to *threshold*
 - 5: **if** $\alpha == 0$ **then**
 - 6: add this node to *isActive*
 - 7: **end if**
 - 8: **end for**
 - 9: *saturation* $\leftarrow 0$
 - 10: *lastlen* \leftarrow length of *isActive*
 - 11: **while** NOT *saturation* **do**
 - 12: *pulse* $\leftarrow 0$
 - 13: **for** each inactive node *v* and *s*'s precursor is active **do**
 - 14: *pulse* \leftarrow *incoming*[*v*]
 - 15: **end for**
 - 16: **if** *pulse* \leq *threshold*[*v*] **then**
 - 17: add node *v* to *isActive*
 - 18: **end if**
 - 19: **if** *lastlen* == length of *isActive* **then**
 - 20: *saturation* $\leftarrow 1$
 - 21: **end if**
 - 22: *lastlen* \leftarrow length of *isActive*
 - 23: **end while**
 - 24: **return** length of *isActive*
-

- **IC**: Calculate influence of seeds set of Independent Cascade model.

Algorithm 3 IC

Input: *graph_cp, S, incoming***Output:** length of *isActive*

```
copy S to isActive and lastActive
balance  $\leftarrow$  0
3: while NOT balance do
    create empty list newActive
    for each active node v's inactive neighbor w do
6:       if incoming[w]  $\leq$  RANDOM NUMBER alpha
           then
               add w to newActive
               add w to isActive
9:       end if
    end for
    if length of newActive == 0 OR length of isActive
        == node_num then
12:       balance  $\leftarrow$  1
    end if
    end while
15: return length of isActive
```

- **hill_greedy:** Generate optimal seed using Hill Greedy Algorithm

Algorithm 4 hill_greedy

Input: *size, model***Output:** *seeds*

```
create empty ans_seeds
cnt  $\leftarrow$  0
while cnt  $\neq$  size do
4:   copy ans_seeds to cur_seed
      high  $\leftarrow$  0
      point  $\leftarrow$  1
      for each node 4 don't be chosen do
8:         add v to cur_seed
           if model == IC then
               new_high  $\leftarrow$  IC(graph_cp, cur_speed, incoming)
           else
12:          new_high  $\leftarrow$  LT(graph_cp, cur_speed, incoming)
           end if
           if new_high  $\leq$  high then
               high  $\leftarrow$  new_high
16:          point  $\leftarrow$  v
           end if
      end for
      add point to ans_seeds
20:   cnt  $\leftarrow$  cnt + 1
end while
return ans_seeds
```

3. Empirical Verification

Empirical verification is confirmed in OJ system. *ISE.py* almost pass all dataset with reasonable bias.

IMP.py only was test with network-5-IC and provide reasonable seeds set.

3.1. Design

SVM classify data at less 90 percent robust.

Random gradient return in reasonable time than simple gradient down algorithm.

Independent cascade model return reasonable influence is small graph and big graph. But Linear Threshold model return less reasonable due with large graph. Finally, I found my code would added some element repeatedly result in larger bias.

For hill greedy, because I just evaluate once, it produce seeds set fast. In small graph, which produce seeds in 65 percent effect. But Not optimal seeds are produced every time.

3.2. Data and data structure

Dictionaries and lists are used widely rather than matrix. Because the input graph is sparse, dictionary are always used to store graph as adjacent list. Global variable *graph_cp*, *graph_cp* and *incoming* are dictionary. Lists always store routes and edges information for local variable.

3.3. Performance

Following table show different performance with different dataset. Offline test perform at Fedora 29 with Intel[®] Xeon(R) CPU E5-1680 v3@3.20GHz and 32GiB memory.

Dataset	Run Time(s)	Result
network-seeds-IC	23.23	5.015
network-seeds-LT	23.23	5.015
network-seeds2-IC	23.24	30.47
network-seeds2-LT	57.79	37.03
NetHEPT-5seeds-LT	58.55	341.9
NetHEPT-5seeds-IC	27.87	276.3
NetHEPT-50seeds-LT		
NetHEPT-50seeds-IC	27.87	1003
network-5-IC	0.73	19.45

3.4. Result

Solutions are accepted most all public dataset' test. But Linear Threshold model produce influence with large bias.

3.5. Analysis

Because natural Hill Greedy is kind of greedy algorithm without heuristics, most of solutions are not optimal. During using adjacent list rather than adjacent matrix, space cost is more lower than $O(n^2)$. Because every point only is evaluate once, only less time it cost. Total, no more than $O(n^2)$ is required.

Acknowledgment

Thanks TA Yao Zhao who explain question and provide general method to solve it. And I also thanks for Kebin Sun discuss and share the idea about algorithm.

References

- [1] Kempe D, Kleinberg J, Tardos É. Maximizing the spread of influence through a social network[C]//Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining. ACM, 2003: 137-146.