

Data Mining Package Report

Kandati Vishnu Sai, Midhilesh E., Sanjay Seetharaman

November 2020

Abstract

With more than 2 billion monthly active users, social networking websites have dramatically changed the lives of people. With a range of features and formats, they cater to the interests of a wide range of users. People build social networks or social relationships with other people who share similar personal or career interests, activities, backgrounds or real-life connections. The availability of services on a wide range of devices such as mobile phones, laptops, and computers has made it more accessible than ever. A user can easily access other users' behaviors and is in turn influenced by them. As a result, effective social influence prediction is vital for a variety of applications such as online recommendation and advertising.

In this report, we discuss DeepInf¹, an end-to-end framework that learns users' latent feature representation for predicting social influence.

¹Qiu, Jiezhong and Tang, Jian and Ma, Hao and Dong, Yuxiao and Wang, Kuansan and Tang, Jie. 2018. "Deepinf: Social influence prediction with deep learning" in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp.2110–2119.

Contents

1	Introduction	1
2	Graph Mining	2
3	DeepInf: Social influence prediction	3
3.1	Introduction	3
3.2	Model framework	3
3.2.1	Sampling Near Neighbour	3
3.2.2	Neural Network	4
3.3	Evaluation Metrics	5

Chapter 1

Introduction

Hi! This is the section for introduction. Create your own tex file and include the corresponding reference in main.tex.

Chapter 2

Graph Mining

Hi! This is a new section.

Chapter 3

DeepInf: Social influence prediction

3.1 Introduction

DeepInf focus on the prediction of user-level social influence. It aims to **predict the action status of a user given the action**. It is a deep learning based framework to represent both influence dynamics and network structures into a latent space and tries to minimize the negative likelihood that was defined in the section 1.

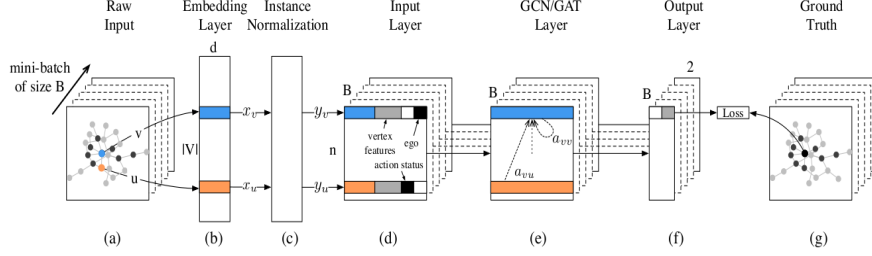
3.2 Model framework

3.2.1 Sampling Near Neighbour

Give a user v , a r-ego network \mathcal{G}_v^r is extracted using breadth-first search (BFS) starting from user v . However, \mathcal{G}_v^r may have different size due to the small world property in the social network. Since most deep learning models expects fixed size data, the graph \mathcal{G}_v^r can be sampled to fixed size.

For sampling a fixed size graph **random walk the restart** (RWR) was used. RWR algorithm is defined as following steps.

- Start random walks from either the ego user v or one of her active neighbors randomly.
- The random walk iteratively travels to its neighborhood with the probability that is proportional to the weight of each edge.
- At each step, the walk is assigned a probability to return to the starting node, that is, either the ego user v or one of v 's active neighbors.



- Run the algorithm until a fixed number of vertices denoted by $\hat{\Gamma}_v^r$ with $|\hat{\Gamma}_v^r| = n$.

After running this algorithm, a sub-graph $\hat{\mathcal{G}}_v^r$ and denote $\hat{S}_v^t = \{s_u^t : u \in \hat{\Gamma}_v^r\}$ be the action statuses of v 's sampled neighbours. Therefore we re-define the optimization objective in section 1 as:

$$\mathcal{L}(\theta) = - \sum_{i=1}^N \log(P_{\theta}(s_{v_i}^{t_i + \Delta t} | \hat{\mathcal{G}}_{v_i}^t, s_{v_i}^{t_i})) \quad (3.1)$$

3.2.2 Neural Network

With the retrieved $\hat{\mathcal{G}}_{v_i}^t$ and \hat{S}_v^t for each user, an effective neural network model to incorporate both the structural properties in $\hat{\mathcal{G}}_{v_i}^t$ and action statuses in \hat{S}_v^t .

Deepinf neural network model consist of network embedding layer, instance normalization layer, input layer and GCN layer as described in the above diagram.

Embedded layer

Network embedding technique encode network structural properties into low dimensional matrix $\mathbf{X} \in \mathcal{R}^{\mathcal{D} \times |V|}$. Deepinf uses Deepwalk algorithm for mapping each users into $\mathcal{R}^{\mathcal{D}}$ space.

Instance Normalization

Instance normalization can remove instance-specific mean and variance, which encourages the downstream model to focus on users' relative positions in latent embedding space rather than their absolute positions. It also prevents overfitting.

Let x_u be the low dimension representation for the user $u \in \hat{\Gamma}_v^r$, the instance normalized vector y_u is obtained by

$$y_{ud} = \frac{x_{ud} - \mu_{ud}}{\sqrt{\sigma_d^2 + \epsilon}} \quad (3.2)$$

for each embedding dimension $d = 1 \dots \mathcal{D}$, where

$$\mu_d = \frac{1}{n} \sum_{u \in \hat{\Gamma}_v^r} x_{ud}, \quad \sigma_d^2 = \frac{1}{n} \sum_{u \in \hat{\Gamma}_v^r} (x_{ud} - \mu_{ud})^2 \quad (3.3)$$

Input Layer

Input later constructs feature vector for each user. Along with the output from the instance Normalization, input layer adds two binary variable. The first variable indicates user's action status and second variable indicates whether the user is the ego user.

GCN

3.3 Evaluation Metrics