

Relationship prediction and link detection of socially connected nodes

Moitrayee Chatterjee
CS Department
Texas Tech University
Lubbock, Texas
Email: moitrayee.chatterjee@ttu.edu

Supra Jyotsna Jampa
Computer Science
Texas Tech University
Lubbock, Texas - 79415
suprajyotsna.jampa@ttu.edu

Abstract—In real world networks of people, it is observed that we rarely tag our contacts. For e.g. if someone is connected to another 5 people, then it is a rare occasion (s)he mentions who among those 5 connections are his/her "school-friend" or "colleagues" or "supervisor" or "family". However, interpreting this type of relationships can be extremely helpful for a number of further applications. For example, if a fund raiser campaign needs to be reached to maximum number of people, and we can detect a 'trust' relationship from the social network of a person, who has donated in that campaign, then the fund raiser can be approached to those connections who has a 'trust' relation with this donor. In this project work, we aim to predict the relationships between users, who have not tagged each other.

Keywords- Link Prediction, Social Ties, Machine Learning, Matrix factorization.

I. INTRODUCTION

In real life there are many social networks like Facebook, Twitter, LinkedIn etc. Which are growing every day. All the ties (links) in online networks are very similar to the real life ties but one main difference is the real life ties are labeled (eg. FRIENDS, FAMILY) while in online networks people do not take time to label them. Many works have focused on mining the relationships in specific domains and one main problem with networks is that there is no labeled data available making it difficult to infer relationships.

In our work, we aim to learn the nature of relationships from overlapping human interactions using mathematical method **Matrix Factorization**.

II. LITERATURE REVIEW

A. Human Mobility, Social Ties, and Link Prediction [1]

In their work, the authors have considered "Spatio-Temporal Co-occurrences" in socially connected users. Their work concentrates on an internet setting and locate those geographic co-events which potentially shapes the derivations about social ties: The fact that two individuals were proximate at only a couple of unmistakable areas at almost similar circumstances can demonstrate a high likelihood that they are socially connected. Utilizing a dataset of 38 million geo-labeled photographs from Flickr, the authors found that the likelihood of a social tie as the number of co-events (k) increases sharply and the transient range (t) diminishes. However, the estimations of the probabilities with respect to the baseline likelihood

of having a social tie is noteworthy. A probabilistic model is used for learning how movement is related to social ties. The authors, used a straightforward model initially and then observed the data closely using a richer one. The initial model has N isolated geographic cells and M individuals, each having one social tie, so that the interpersonal organization comprises of (M^2) disjoint edges. Every day, each combination of companions visits a place mutually with likelihood β and independently with likelihood $1-\beta$; in either case the decision of location(s) is made arbitrarily. Utilizing Bayes' Law, the likelihood that two individuals are companions (occasion θ) given that they visit the very same cells on k back to back days (occasion D) is:

$$P(\theta|\mathbf{D}) = P(\theta) \frac{P(\mathbf{D}|\theta)}{P(\mathbf{D})} \quad (1)$$

The prior probability that two individuals are companions, $P(\theta)$ is $1/(M-1)$, while the likelihood of two connected individuals being at a similar place on a given day is calculated as:

$$p = (\beta) + \frac{1-\beta}{N} \quad (2)$$

This basic probability model captures the macro features but fails to take into account the micro features. Hence the authors take "homophily" into account. This is the way individuals associated in social ties will probably take part in related exercises, because of their inalienable closeness. The N geographic cells are organized in a framework, and each match of companions (A, B) has a randomized "home" cell, drawn from the 2D experimental dispersion of Flickr photos. Results of this model recommend an approach to evaluate such overlap.

B. Inferring social ties from geographic coincidences [2]

The authors of this paper studied the extent to which the mobility pattern of individual influence the social ties. They used the **Call Detail Report** of an anonymous country. The CDR had data for 3 months of activity of 6million users. If there is a pair of user u and v , then in their work the author addressed three questions: 1. "How similar is the movement of u and v ", 2. "How connected are u and v in the social network", 3. "How intense is the interaction between u and v ".

Network proximity and Mobile homophily is used in this work for link prediction.

Network proximity: Proximity in mobile network can be used for the future link prediction between two user. From the location obtained from mobile data, if two users are geographically closely located, then there a possibility of them to be connected in near future. The proximity in network is studied using four different measures:

- Common neighbors: Common users who are connected to both u and v .
- Adamic-Adar: Instead of only counting the common connection, each connected person is given a weight.
- Jaccards coefficient: It is the size of the intersection of connections of u and v .
- Katz: Summation over every single conceivable way from u to v with "exponential damping by length to weight short paths".

Mobile Homophily: In graph based methodologies, predicting new connections by utilizing mobility data can be achieved by searching the degree of closeness in physical space between two people. However, in reality, individuals who share high degree of overlap in their connections, are expected to have a better likelihood of forming new links. The authors used following measures to define similar mobility pattern between individuals:

- Distance.
- Spatial Co-Location Rate. The likelihood of u and v to visit the same location, but may not be at same time.
- Spatial Cosine Similarity: The cosine similarity of user x and y s trajectories, capturing how similar their visitation frequencies are, assigned by the cosine of the angle between the two vectors of number of visits at each location for x and y .
- Weighted Spatial Cosine Similarity. The tf-idf version of cosine similarity of the visitation frequencies of users x and y , where the contribution of each location l is inversely proportional to the (log of) its overall population in l .
- Co-Location Rate: Both spatial and temporal vicinity is taken into account and normalized by the number of times u and v were in same time frame.
- Weighted Co-Location Rate: The likelihood of u and v to co-exists in same time at same location and it is normalized by the population density of that co-location at the given time.
- Extra-role Co-Location Rate: The likelihood of u and v to co-locate in the same hour of night or over weekends.

The association between individual mobility patterns and social vicinity in the call graph is studied by measuring the connection between u and v utilizing aforementioned procedures. The quality of the ties in the network is measured by the numbers of calls set between any two individual and also taken into account. Finally, they have formalized the link prediction as a binary classification problem where the network and mobility parameters are used as predictive variables. In short,

connections observed in the "past" are used to predict the "future" links.

C. Inferring social ties in Heterogeneous Networks [5]

This paper focuses on a systematic investigation is conducted on the different networks to define problem and they proposed a Transfer based factor graph (TranFG) model. The model is a framework that considers social theories and applies it to infer social ties in target network.

- Factors influencing the social ties in a network: To infer social ties in a domain specific networks we can devise methods for each networks which cannot be used in heterogeneous networks. Few statistics according to social psychological social theories are used to analyze the network:

Social Balance: In a network social balance theory can be applied to a triad of nodes, the balanced when there is only one pair of friends or all three of them are friends.

Structural Holes: A person in a network is said to span a structural hole if he/she is connected to many other nodes in the network which are not so well connected. This is a very important feature of a node as a node which is a structural hole can be used this node to spread information to its connected nodes. Generally structural nodes are identified for better promotion. eBay is a website which is a structural hole which is connected to many people who may not be well connected to one other. The main is to test if a structural hole tends to have the same type of relationship with the other users. (i)For a node, the total number of pairs of neighbors which are not connected directly are calculated. (ii)The nodes are ranked based on their values and the top 1 percent are considered as structural holes. (iii)If a nodes is a structural hole are the users or most likely 70 percent of the users have same relationship with the node. (iv)Mostly disconnected users have similar type of relationship with the structural hole.

Social Status: In this theory we consider directed graphs and each link is given a +/- sign based on where the target node has higher/lower social status than the source node. When we consider a triangle and each negative sign edge is reversed and the sign the resulting triangle is acyclic.

Opinion Leader: A opinion leader is node in the network which gets the information first and then passed on to other nodes that are connected to it. These nodes are the ones that are in the higher nodes in a hierarchical networks like networks of employees in a company. (i)The users are categorized into opinion leaders and opinion users by PageRank (ii)The users are ordered and the 1 percent users having highest Page Rank scores are considered as opinion leaders Social theories into semi supervised learning framework is used to infer the social ties in target network.

- Methodology: Two graphs: Source Network and target network which have few labeled and few unlabeled edges-undirected graph is considered and later a solution for

directed graph is proposed. The model is such formulated such that a function $f : (GT—GS) \rightarrow YT$

To infer the relation types in the target network using the supervised information (ie. The labeled of the links in the source network). Input: Two partially labeled networks such that the labels in the source network is much greater than the target network. The type of the relation in the two networks might be very different, even the labels. Since the networks are not completely labeled the predictive function should consider the unlabeled information also. To infer social ties five different networks are considered,

- Epinions: A network of reviewers where the nodes are users who review products and rate other users reviews with a trust or distrust which form the links.
- Slashdot: Slashdot is a website where users(nodes) share technology news and can tag other people as friends or foes(links).
- Mobile: A network of mobile users who are connected with each through calls or texts in the same place. The users are the nodes and there is a link if they communicate with each other.
- Coauthor: A network from Arnetminer.org which consists of authors and co- authors.
- Enron: A network of communication in an hierarchical organization of Enron where the employees communicate with each other through emails and the links are of type manager and subordinate.

Observations when the social theories are applied on the graphs : Domain specific features will change in the case of multiple heterogeneous networks The problem is connected to psychological theories the analysis is conducted on the network based correlations of the following statistics like how social balance and social status are satisfied, the behavior of structural in different networks and the effect of opinion leaders on other networks are used to analyze the network. MODEL:A network is represented as $[G = V, E^L, E^U, X]$

Where, E^L is the set of labeled edges and E^U is the set of unlabeled edges and $E^L + E^U = E$, the total number of edges. X is the attribute matrix. The formulation is:

$$P(Y|X, G) = \frac{P(X, G|Y)P(Y)}{P(X, G)} \quad (3)$$

When we assume that probability of attributes when the each link's label is conditionally independent, according to Bayes Rule we have

$$P(Y|X, G) = \frac{P(X, G|Y)P(Y)}{P(X, G)} \propto P(X|Y).P(Y|G) \quad (4)$$

To instantiate the probabilities we use Hammersley-Clifford theorem. Using which we get α_j = weight of the jth attribute μ_k = weight of the kth correlation feature function

For the network G with labeled information Y , learning the predictive function is to estimate a parameter configuration,

$$\odot(\theta) = \log P_\theta(Y|X, G) \quad (5)$$

- Model learning and Inferring: To estimate the parameter configuration and to maximize the objective function. A gradient decent method (Newton-Raphson) is used to solve the objective function. An algorithm is devised based on which the networks are evaluated.

- Algorithm:

Input: a source network G_S , a target network G_T , and the learning rate η

Output: estimated parameters $\theta = (\{\alpha\}, \{\beta\}, \{\mu\})$

Initialize $\theta \leftarrow 0$;

Perform statistics according to social theories;

Construct social theories based features $h_k(Y_c)$;

repeat

Step 1: Perform LBP to calculate marginal distribution of unknown variables in the source network $P(y_i|x_i, G_S)$;

Step 2: Perform LBP to calculate marginal distribution of unknown variables in the target network $P(y_i|x_i, G_T)$;

Step 3: Perform LBP to calculate the marginal distribution of clique c , i.e., $P(y_c|X_c^S, X_c^T, G_S, G_T)$;

Step 4: Calculate the gradient of μ_k according to Eq. 8 (for α_j and β_j with a similar formula);

Step 5: Update parameter θ with the learning rate η :

$$\theta_{\text{new}} = \theta_{\text{old}} + \eta \cdot \frac{\odot(\theta)}{\theta}$$

until Convergence;

- Experimental Setup: For the five types of networks, SVM, CRF, PFG, TranFG methods are used for comparison. The experiments conducted for the same features and for different pairs of (source and target) networks, Precision, Recall and F1- Measure are used for evaluation.

Data Set	Method	Prec.	Rec.	F1-score
Epinions (S) to Slashdot (T) (40%)	SVM	0.7157	0.9733	0.8249
	CRF	0.8919	0.6710	0.7658
	PFG	0.9300	0.6436	0.7607
	TranFG	0.9414	0.9446	0.9430
Slashdot (S) to Epinions (T) (40%)	SVM	0.9132	0.9925	0.9512
	CRF	0.8923	0.9911	0.9393
	PFG	0.9954	0.9787	0.9870
	TranFG	0.9954	0.9787	0.9870
Epinions (S) to Mobile (T) (40%)	SVM	0.8983	0.5955	0.7162
	CRF	0.9455	0.5417	0.6887
	PFG	1.0000	0.5924	0.7440
	TranFG	0.8239	0.8344	0.8291
Slashdot (S) to Mobile (T) (40%)	SVM	0.8983	0.5955	0.7162
	CRF	0.9455	0.5417	0.6887
	PFG	1.0000	0.5924	0.7440
	TranFG	0.7258	0.8599	0.7872

D. Inferring real world connections and discovering, labeling, and characterizing communities[6]

Information has become one of the highly value commodities in the present world. It is used for multiple operations such as optimization of any system, understand the trends and understanding the thought process of people as well. This paper tries to analyze the social interactions between people using the information gathered and also the factors which are the underlying reasons for this structure. This is later used to infer the real world connections and the communities of people on the whole. The paper focuses on understanding these social

interactions and accordingly the inferences through individual homepages. Individual homepages are something which show us the interests and experiences of human beings. While it is difficult to find individual homepages as they are not free floating on the web, they can be used to infer regarding the communities on the whole. These pages also tend to consist links which consequently link them with people of similar interests and basically form communities of people. This inference can be considered as due to side effect of studying the homepages of people. This is in turn a side effect of people sharing their information on the digital world. In the coming sections, we will discuss about the structure of webpages in terms of small world phenomena, later visualizing the interface for exploring the networks between people.

The short world phenomena speaks about how people in the world are connected through a small chain of acquaintances. It basically means that people with similar interests tend to spend time together in smaller circles. Extrapolating from this point of view, it can also be said that people on world wide web and social network on the whole, belong to a small world of its own. In this paper, research has been conducted on homepages of MIT and Stanford students to analyze their networks, interests and accordingly make inference regarding the communities mentioned before. From the study conducted, it was noticed that 70 percent of the MIT users and 30 percent of the Stanford users are connected to other. While people with similar interests might be linked, it is also possible that people with similar links are not connected either. This was successfully eliminated from the scope of study as well. It was also noticed that while users provided out links to two or three, there were users who received many out links. Similarly, there were few users who had multiple links as compared to one or two for most of the users. Clustering coefficient was used to measure the extent to which users band together. It showed that the results for MIT and Stanford was fairly similar whereas the geodesic distance was more for Stanford as compared to MIT. We can infer from these results that MIT and Stanford are part of small world networks where they are closely linked based on the geodesic distance.

After creating the networks, in order to understand the reason behind the links, an interface was created which helps them search for the users, notice their details, visualize their network, match the users to others based on their links, text and mailing lists. Once the interface is created, relations can be predicted using a matchmaking algorithm created. For example, likelihood of link between users can be done by finding things in common and also by assigning higher weights to things which are pertained to just certain people.

III. DATASET DESCRIPTION

We are using **Slashdot social network, February 2009** [3] dataset. Slashdot is a website for sharing technology related news. On Slashdot one user submit one technology-related news and editor evaluates them. The Slashdot Zoo features allows the user to tag each other as *friends* and *foes*. In our project we aim to infer relationships between users who have

not tagged relationship with each others. The dataset is a directed graph of 82168 nodes representing users and 948464 edges representing *friendship* relation.

IV. TECHNICAL APPROACH

A wide variety of algorithms are developed to aid the learning process. However, matrix factorization is an effective mathematical approach that can help in learning the latent features of human ties. Tensor factorization could also be used for solving our problem, however, it would have more complex and would need better understanding of the implementation and working principles.

A. Basic Idea

The term "Matrix factorization" refers to the fact of finding two matrices out of one, in a way such that when those two matrices are multiplied, we get back the initial matrix. NetFlix uses matrix factorization for their movie recommendation system.

B. Mathematics of Matrix Factorization

We consider our dataset of Slashdot users as a boolean matrix where both columns and rows represents the nodes or users. Each element in the matrix is 1 if two users/nodes have a friendship edge between them and 0 in case of no connection. So, we have to find two matrices **P** and **Q** such that their product will approximate the original matrix **R**. Mathematically:

$$\mathbf{R} \approx \mathbf{P} \times \mathbf{Q}^T = \hat{\mathbf{R}}$$

C. Implementation

We will implement the matrix factorization algorithm in Python using the **Scikit-learn** [4] and the standard 0.01 learning rate will be used. However, Depending upon weight and biases this learning rate may vary. At the end of all the iterations of the algorithm a new boolean matrix will be obtained and some of the intersection elements containing 0's will now contain 1, signifying possible *friendship*.

V. FUTURE WORK

Currently the code development is in progress. In parallel the hyper parameters are also in preparation phase. 0.8 fraction of the total dataset will be used for training and rest 0.2 will be used as test dataset.

ACKNOWLEDGMENT

We would like to thank our professor Dr. Mahshid R. Naeini for guiding us through out our project.

REFERENCES

- [1] Wang, Dashun, et al. "Human mobility, social ties, and link prediction." Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining. ACM, 2011.
- [2] Crandall, David J., et al. "Inferring social ties from geographic coincidences." Proceedings of the National Academy of Sciences 107.52 (2010): 22436-22441.
- [3] J. Leskovec, K. Lang, A. Dasgupta, M. Mahoney. Community Structure in Large Networks: Natural Cluster Sizes and the Absence of Large Well-Defined Clusters. Internet Mathematics 6(1) 29–123, 2009.

- [4] Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Pret-tenhofer, R. Weiss, and V. Dubourg. scikit-learn: Machine learning in Python. *The Journal of Machine Learning Research*, 12:28252830, 2011.
- [5] J. Tang, T.Lou, J. Kleinberg "Inferring Social Ties across Heterogenous Networks".
- [6] L. A. Adamic and E. Adar. Friends and neighbors on the web. *SOCIAL NETWORKS*, 25:211230, 2001.