Welcome to

DS504/CS586: Big Data Analytics Graph Mining II Prof. Yanhua Li

Time: 6:00pm –8:50pm Mon. and Wed. Location: SL105 Spring 2016

Reading assignments

We will increase the bar a little bit Please add more of your ideas, and share them with us in the class.

Project 1

Team 1:

Utilizing the data from Allstate to predict the possible Bodily Injury Liability Insurance claim payments that the company may pay in line with the vehicle features.

Team 2:

3D Video Growing Trend on YouTube

Project 1

Team 3:

Designing a sampling method that can estimate the available capacity of hotels in a large area and a certain time range as accurate as possible.

Team 4:

A huge number of professional as well as amateur programmers use Stackoverflow.com to find solutions to their questions regularly. Thus, our team seeks to give such users a general idea of which skills to learn to create their projects and how to manage their skills by offering a projectrelated heat map of programming languages and knowledge points.

Project 1

Team 5:

MEASURING RESTAURANT DIVERSITY INDEX FOR DIFFERENT CITIES in Yelp

Team 6:

We will estimate the number of valid and invalid users among the population. Lastly, we will analyze other related web site statistics such as passive and active users.



Real-life graph contains complex contents – labels associated with nodes, edges and graphs.



Node Labels:

Location, Gender, Charts, Library, Events, Groups, Journal, Tags, Age, Tracks.

Large Scale Graphs.

Facebook Twitter LinkedIn Last.FM LiveJournal del.icio.us DBLP

of Users 400 Million 105 Million 60 Million 40 Million 25 Million 5.3 Million 0.7 Million

of Links 52K Million **IOK Million** 0.9K Million **2K Million 2K Million** 0.7K Million 8 Million

Mining in Big Graphs

Network Statistic Analysis (last lecture)

- Network Size
- Degree distribution.

- Node Ranking (this lecture)
 - Identifying most important/influential nodes
 - Viral Marketing, resource allocation

Characterize Node Importance

- Rank the webpages in search engine.
- Viral Marketing, resource allocation
- Open a new restaurant, find the optimal location



Ranking nodes on an undirected graph



They are equivalent.

Ranking nodes on a directed graph



They are equivalent?

Random Walk (Undirected Graph) Adjacency matrix $A = \begin{pmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \end{pmatrix} D = \begin{pmatrix} 3 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 2 \end{pmatrix}$ Symmetric 4 Transition Probability Matrix Undirected $P = A \bullet D^{-1} = \begin{pmatrix} 0 & 1/3 & 1/3 & 1/3 \\ 1/2 & 0 & 1/2 & 0 \\ 1/3 & 1/3 & 0 & 1/3 \\ 1/2 & 0 & 1/2 & 0 \end{pmatrix}$ $P_{ij} = \frac{1}{k_i}$ $x_{t,i} = \sum_{j}^{i} x_{t-1,j} p_{ji}$ ✤ |E|: number of links π(1)=3/10 Stationary Distribution π(3)=3/10 π(2)=2/10 $\pi_i = \frac{a_i}{2|E|}$ π(4)=2/10

Random Walk (directed graph) Strongly Connected Graphs & Aperiodic * Adjacency matrix

$$A = \begin{pmatrix} 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 1 & 0 & 0 & 0 \end{pmatrix} D = \begin{pmatrix} 2 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$
Asymmetric

Transition Probability Matrix

$$P_{ij} = \frac{1}{k_{out,i}}$$

$$x_{t,i} = \sum_{j} x_{t-1,j} p_{ji}$$

IE: number of directed links

Stationary Distribution

$$\pi_i \neq \frac{d_i}{2|E|}$$



$$P = A \bullet D^{-1} = \begin{pmatrix} 0 & 1/2 & 1/2 & 0 \\ 1/2 & 0 & 1/2 & 0 \\ 1/3 & 1/3 & 0 & 1/3 \\ 1 & 0 & 0 & 0 \end{pmatrix}$$

 $\pi(1)=6/18=1/3$ π(2)=4/18=2/9 $\pi(3)=3/18=1/6$ π(4)=5/18

Ranking nodes in a directed graph



They are no longer equivalent.

directed graphs Strongly Connected Graphs & Aperiodic

- Periodic
- VS
- Aperiodic Graphs
 - The greatest common divisor of the lengths of its cycles is one or not
- Disconnected graph
- VS
- Connected graph
 - Strongly Connected
 - VS
 - Weakly Connected







🚼 Everything

💷 Images

Videos

Shopping

Durham, NC

🕅 News

More

Dan Teaque

About 999,000 results (0.12 seconds)

Search Advanced search

🕨 Dr. Dan Teague 🔍

Dan Teague has been an insitucior of Maihematics at the North Carolina School of Science and Mainematics since 1982. He received his undergraduale degree ... www.maa.org/abouimaa/feague_bio.thimi - Cached

NCSSM Math Department Taks and Papers 🔍

by Dan Teague, 2010. Gas Station Problem by Dan Teague, 2008. The Cookie Problem ... by Dan Teague, Markov, Chebyshev, and the Weak Law of Large Numbers.

www.ncs.sm.edu/courses/maih/Taiks/Index.him - Cached - Similar

Dan Teague's Page 🦳

Dan Teague, I am an insitucior of Mathematics at the North Carolina School ... courses incost ineduitrain/Faculty/D_Teague Jrimi - Cached - Similar

Show search loois

Change location

HiShow more results from nossmiedu

Big Dan Teague (Cikaracter) 🦳

Big Dan Teague (Characler) from O Brother, Where Arl Thou? (2000), Share his page : ... (2000) Big Dan Teague : So long boys . See you in the funny papers www.imdb.com/characler/ch0004822/ - Cached - Similar

O Brother, Where Art Thou? - Wikibedia, the free encyclopedia 🤐

John Goodman as Daniel "Big Dan" Teague, one of the main enemies of the Irlo in he film. Masquerading as a Bible salesman, he cons Everell, hen robs him.... en.wikipedia.org/wiki/0_Brother_Where_Art_Thou%3F - Cached - Similar

Dan Teague profiles | Linkedin 🤐

View the profiles of professionals named Dan Teague on Linked in. There are 17 professionals named Dan Teague, who use Linked in to exchange information, ... www.linkedin.com/pub/dir/Dan/Teague - Cached

Dan Teague - Linkedin 🔍

Louisville, Kenlucky Area - Purchasing Manager al Beam Global Spirils & Wine View Dan Teague's professional profile on Unkedin. Unkedin is the world's www.linkedin.com/pub/dan-feague/8/496/124

Show more results from linkedim.com

Dan Teague | Facebook 🔍

Dan Teague is on Facebook, Join Facebook to connect with Dan Teague and others you may know. Facebook gives people the power to share and makes the world ... www.facebook.com/people/Dan-Teague/1511802817 - Cached

Dan Teague (DTeaglie 4593) on Twitter 🔍

Gelishori, Imely messages from Dan Teague. Twiller is a rich source of instanly updaled information. It's easy to slay updaled on an incredibly wide ... willer.com/DTeague 4593 - Cached

Videos for Dan Teague





5 mln - Nov 17, 2010 Uploaded by KeyCurriculum Press youlube.com

Why This **Order**?



They are no longer equivalent.

Naïve PageRank

Adjacency matrix

Transition Probability Matrix

$$P_{ij} = \frac{1}{k_{out,i}}$$

$$R_i = \sum_{j} R_j p_{ji}$$

Stationary Distribution

$$P = A \cdot D^{-1} = \begin{pmatrix} 0 & 1/2 & 1/2 & 0 \\ 1/2 & 0 & 1/2 & 0 \\ 1/2 & 0 & 1/2 & 0 \\ 1/3 & 1/3 & 0 & 1/3 \\ 1 & 0 & 0 & 0 \end{pmatrix}$$

$$\pi(1) = 6/18 = 1/3$$

$$\pi(2) = 4/18 = 2/9$$

$$\pi(3) = 3/18 = 1/6$$

$$\pi(4) = 5/18$$

 $R_i = \pi_i$

Disconnected Graph & Random surfing behaviors

Standard PageRank

Adjacency matrix

$$A = \begin{pmatrix} 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 1 & 0 & 0 & 0 \end{pmatrix} \quad D = \begin{pmatrix} 2 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

Transition Probability Matrix (d=0.85)

$$P_{ij} = \frac{1}{k_{out,i}}$$

$$P = A \cdot D^{-1} = \begin{pmatrix} 0 & 1/2 & 1/2 & 0 \\ 1/2 & 0 & 1/2 & 0 \\ 1/3 & 1/3 & 0 & 1/3 \\ 1 & 0 & 0 & 0 \end{pmatrix}$$

$$R_i = d \sum_j R_j p_{ji} + (1 - d) \frac{1}{n}$$

0.4625

0.0375

0.0375

0.3208 0.0375

0.0375

0.0375

0.3208

0.8875

0.4625

0.0375

0.0375

0.0375

0.8875

0.3208

0.0375

Stationary Distribution (J is all-1 matrix).

$$R_i = \pi_{pr,i} \qquad \qquad P_{pr} = d \cdot P + (1-d)\frac{1}{n}J =$$

- Convergence
 - Leading eigenvector of P_{pr}

How to quantify the importance as a hub and authority separately?

Hub & Authority (HITS)

Adjacency matrix

- Hub and authority
 - Initial Step: hub(p) = 1; auth(p) = 1;
 - Each step with normalization:

$$hub(p) = \sum_{i=1}^{n} auth(i); \qquad hub(p) = \frac{hub(p)}{\sqrt{\sum_{i=1}^{n} hub(i)^{2}}};$$

$$auth(p) = \sum_{i=1}^{n} hub(i); \quad auth(p) = \frac{auth(p)}{\sqrt{\sum_{i=1}^{n} auth(i)^2}};$$

- Convergence
 - hub and authority are the left and right singular vector of the adjacency matrix A.



A Note on Maximizing the Spread of Influence in Social Networks

E. Even-Dar and A. Shapira





Social Influence



Social Influence



A few days ago, I had a chance to try a particular brand of frozen dessert maker, which was only available in the United States. I sent out a "tweet" in China on Sina Weibo, sharing

To my surprise, this tweet was retweeted over 100,000 times, as Chinese young people and parents showed great interest in having such a machine. Even more surprising, within 3 hours, Taobao (China's eBay) had hundreds of sellers, offering to buy such a machine in the US and shipping it to China for the buyer. Even though the shipping cost was more than the machine, thousands were sold within a day (one store reported 51 sales, and there were over

Voter Influence Model

Opinion diffusions

Switch opinions back and forth

Word of mouth effect!

Randomly selecting one neighbor to adopt its opinion



[1] P. Clifford and A. Sudbury. A model for spatial conflict. *Biometrika*, 60(3):581, 1973.

Influence Maximization

Budget: Selecting k individuals as initial red seeds Assumption: Uniform cost of selecting each initial seed Goal: Maximize the number of future red nodes

[15] E. Even-Dar and A. Shapira. A note on maximizing the spread of influence in social networks. In WINE, 2007.

Formulation



Formulation (Random Walk)

Influence at step t:

Influence contribution:Short term $\max_{x_0} : \mathbf{1} x_t^T$ Long term $\max_{x_0} : \lim_{t \to \infty} \mathbf{1} x_t^T$

 X_t^T is a column vector, which is the transpose of row vector X_t

Matrix form:

$$x_t = x_0 P^t$$

$$\lim_{t \to \infty} x_t = \lim_{t \to \infty} x_0 P^t = \pi$$

Influence contribution:

Short term

$$\max_{x_0} : f_t(x_0) - f_0(x_0)$$
$$\max_{x_0} : \lim_{t \to \infty} f_t(x_0) - f_0(x_0) = x_0 \pi^T - f_0(x_0)$$

Long term

 $\mathbf{1}\mathbf{x}_{t}^{T}$

Influence Maximization

Budget: Selecting C for initial red seeds

Assumption: Heterogeneous costs of selecting different initial seeds (c_i) Goal: Maximize the number of future red nodes



[15] E. Even-Dar and A. Shapira. A note on maximizing the spread of influence in social networks. In WINE, 2007.

Knapsack problem

Knapsack problem

Weight = Influence Value/ Stationary distribution Size = Cost c_i of choosing a node n_i



What if directed social graph

One-way connection

Randomly select an out-going neighbor



Adopt the opinion of one of the outgoing neighbors.



Adopt the opposite opinion of foe, the same opinion of friend

[WSDM'13] Yanhua Li, Wei Chen, Yajun Wang, Zhi-Li Zhang, **Influence Diffusion Dynamics and Influence Maximization in Social Networks with Friend and Foe Relationships**. *The 6th ACM International Conference on Web Search and Data Mining*, February 4-8, 2013, Rome, Italy.

Any Comments & Critiques?

Next Class: Graph Mining (Presentation)

- Do assigned readings before class
- Submit reviews/critiques
- Attend in-class discussions