

CarnegieMellon

15-826: Multimedia Databases and Data Mining

Lecture #28: Graph mining -
Belief propagation
Christos Faloutsos

1

CarnegieMellon

Must-read material


- Danai Koutra, Tai-You Ke, U. Kang, Duen Horng Chau, Hsing-Kuo Kenneth Pao, and Christos Faloutsos. [Unifying guilt-by-association approaches: theorems and fast algorithms.](#) ECML/PKDD'11, Athens, Greece. (new)

15-826 Copyright (c) 2019 C. Faloutsos #2

2

CarnegieMellon

Main outline




- Introduction
- Indexing
- Mining
 - Graphs – patterns
 - ➔ – Graphs – generators and tools
 - Association rules
 - ...

15-826 Copyright (c) 2019 C. Faloutsos 3

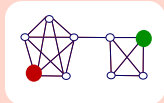
3

CarnegieMellon

Detailed outline




- Graphs – generators
- Graphs – tools
 - Community detection / graph partitioning
 - ➔ – ‘Belief Propagation’ & fraud detection
 - Motivation - Basics
 - Fast, linear approximation (FaBP)
 - Latest: zooBP
 - Success stories
 - Influence/virus propagation & immunization



15-826 Copyright (c) 2019 C. Faloutsos 4


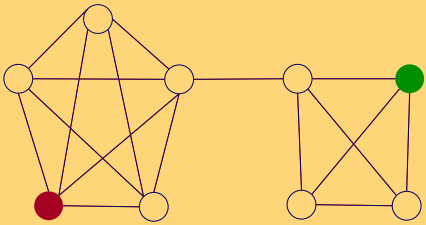
4

CarnegieMellon



Problem


- What color, for the rest?
 - Given homophily (/heterophily etc)?

15-826 Copyright (c) 2019 C. Faloutsos 5


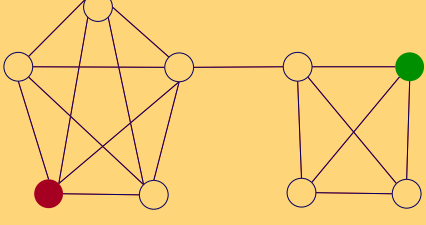
5

CarnegieMellon

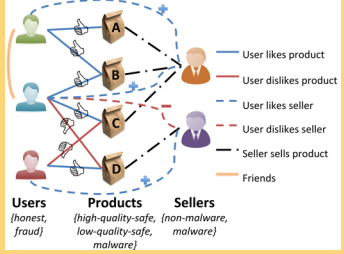


Short answer:

- What color, for the rest?
- A: Belief Propagation ('zooBP')

www.cs.cmu.edu/~deswaran/code/zoobp.zip



Users: {honest, dishonest, fraud}
 Products: {high-quality-safe, low-quality-safe, malware}
 Sellers: {non-malware, malware}


Legend:
 - User likes product (blue solid line)
 - User dislikes product (red solid line)
 - User likes seller (blue dashed line)
 - User dislikes seller (red dashed line)
 - Seller sells product (black solid line)
 - Friends (orange solid line)

15-826 Copyright (c) 2019 C. Faloutsos 6

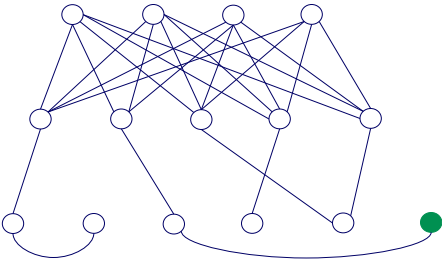
6

CarnegieMellon

E-bay Fraud detection



w/ Polo Chau & Shashank Pandit, CMU
[PKDD'06][WWW'07]

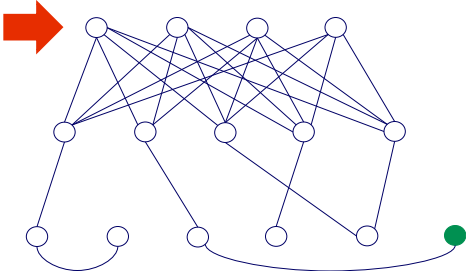


Detecting Fraudulent Personalities in Networks of Online Auctioneers. *Duen Horng (Polo) Chau, Shashank Pandit, and Christos Faloutsos. (PKDD) 2006*

7

CarnegieMellon

E-bay Fraud detection



15-826 Copyright (c) 2019 C. Faloutsos 8

8

CarnegieMellon

E-bay Fraud detection

15-826 Copyright (c) 2019 C. Faloutsos 9


9

CarnegieMellon

E-bay Fraud detection - NetProbe

15-826 Copyright (c) 2019 C. Faloutsos 10

10

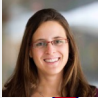


Prof. Danai Koutra
U. Michigan

Background

15-826 Copyright (c) 2019 C. Faloutsos 11

11



Belief Propagation

- Iterative message-based method
- “Propagation matrix”:
 - ◊ Homophily

class of sender	class of receiver	0.9	0.1
0.1	0.9		

until stop criterion fulfilled

AI

PL

[Pearl '82][Yedidia+ '02] ... [Gonzalez+ '09][Chechetka+ '10]

15-826 Copyright (c) 2019 C. Faloutsos 12

12

Ca

Belief Propagation Equations

message($i \rightarrow j$) \approx belief(i) \square homophily strength

●	●
●	●
0.9	0.1
0.2	0.8

15-826 Copyright (c) 2019 C. Faloutsos 13

13

Ca

Belief Propagation Equations

belief of i \cdot $b_i(x_i) \propto \phi_i(x_i) \cdot \prod_{j \in N(i)} m_{ij}(x_i)$

prior belief

messages from neighbors

15-826 Copyright (c) 2019 C. Faloutsos 14

14

Ca




Background

15-826 Copyright (c) 2019 C. Faloutsos 15

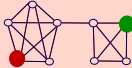
15

CarnegieMellon

Roadmap



- Introduction – Motivation
- Part#1: Graphs
 - ...
 - P1.5: belief propagation
 - Basics
 - • Fast, linear approximation (FaBP)
 - Latest: zooBP



15-826 Copyright (c) 2019 C. Faloutsos 16

16

CarnegieMellon

Unifying Guilt-by-Association Approaches: Theorems and Fast Algorithms



Danai Koutra
U Kang
Hsing-Kuo Kenneth Pao

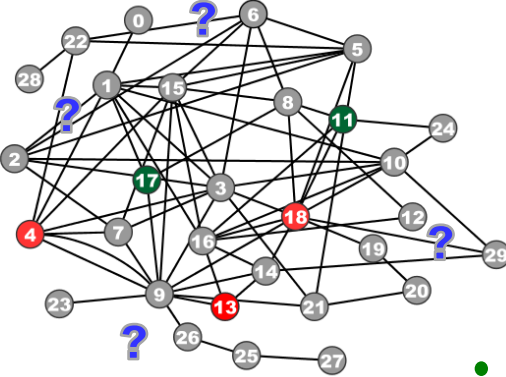
Tai-You Ke
Duen Horng (Polo) Chau
Christos Faloutsos

ECML PKDD, 5-9 September 2011, Athens, Greece

17

CarnegieMellon

Problem Definition: GBA techniques



Given: Graph; &
few labeled nodes

Find: labels of rest
(assuming network effects)

●
●

●
●

15-826

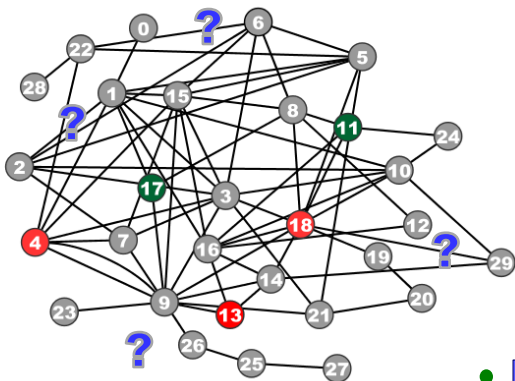
Copyright (c) 2019 C. Faloutsos

18

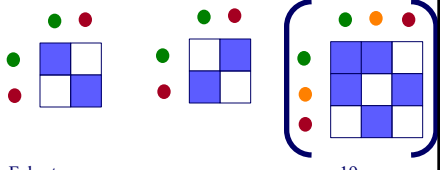
18

CarnegieMellon

Problem Definition: GBA techniques



Given: Graph; & few labeled nodes
Find: labels of rest (assuming network effects)



15-826 Copyright (c) 2019 C. Faloutsos 19

19

CarnegieMellon

BP vs. Linearized BP DETAILS

Original [Yedidia+]:

Belief Propagation

$m_{ij}(x_j) \leftarrow \sum_{x_i} \phi_i(x_i) \cdot \psi_{ij}(x_i, x_j) \cdot \prod_{n \in N(i) \setminus j} m_{ni}(x_i)$

$b_i(x_i) \leftarrow \phi_i(x_i) \cdot \prod_{j \in N(i)} m_{ij}(x_j)$

↑
non-linear

- Closed-form formula?
- Convergence?

19 C. Faloutsos 20

20

CarnegieMellon

BP vs. Linearized BP

DETAILS

Original [Yedidia+]:

Belief Propagation

$$m_{ij}(x_j) \leftarrow \sum_{x_i} \phi_i(x_i) \cdot \psi_{ij}(x_i, x_j) \cdot \prod_{n \in N(i) \setminus j} m_{ni}(x_i)$$

$$b_i(x_i) \leftarrow \phi_i(x_i) \cdot \prod_{j \in N(i)} m_{ij}(x_i)$$

non-linear

✓
Closed-form formula?

✓
Convergence?

Our proposal:

Linearized BP

BP is approximated by

$$[\mathbf{I} + a\mathbf{D} - c'\mathbf{A}] \mathbf{b}_h = \phi_h$$

1		
	1	
		1

d1		
	d2	
		d3

0	1	0
1	0	1
0	1	0

?

0
-10 ⁻²
10 ⁻²

linear

19 C. Faloutsos
21

21

CarnegieMellon

Are they related?


- RWR (Random Walk with Restarts)
 - google’s pageRank (*‘if my parents are important, I’m important, too’*)
- SSL (Semi-supervised learning)
 - minimize the differences among neighbors
- BP (Belief propagation)
 - send messages to neighbors, on what you believe about them

15-826
Copyright (c) 2019 C. Faloutsos
22

22

CarnegieMellon

Are they related? YES!



- RWR (Random Walk with Restarts)
 - google’s pageRank (*‘if my parents are important, I’m important, too’*)
- SSL (Semi-supervised learning)
 - minimize the differences among neighbors
- BP (Belief propagation)
 - send messages to neighbors, on what you believe about them

15-826 23
 Copyright (c) 2019 C. Faloutsos

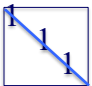
23

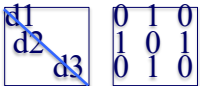
CarnegieMellon

Correspondence of Methods DETAILS


$\mathbf{p} = (1-c)/n [\mathbf{I} - c \mathbf{B}]^{-1} \vec{e}$

Method	Matrix	Unknown	known
RWR	$[\mathbf{I} - c \underline{\mathbf{A}}\mathbf{D}^{-1}]$	\mathbf{x}	$(1-c)\mathbf{y}$
SSL	$[\mathbf{I} + a(\mathbf{D} - \underline{\mathbf{A}})]$	\mathbf{x}	\mathbf{y}
FABP	$[\mathbf{I} + a \mathbf{D} - c' \underline{\mathbf{A}}]$	\mathbf{b}_h	Φ_h






adjacency matrix



final labels/beliefs



prior labels/beliefs

15-826 24
 Copyright (c) 2019 C. Faloutsos

24

CarnegieMellon

Problem: e-commerce ratings

fraud

- User likes product
- User dislikes product
- User likes seller
- User dislikes seller
- Seller sells product
- Friends

Users
(honest, fraud)

Products
(high-quality-safe, low-quality-safe, non-malware, malware)

Sellers
(non-malware, malware)

- Given a heterogeneous graph on users, products, sellers and positive/negative ratings with “seed labels”
- Find the top k most fraudulent users, products and sellers

15-826 Copyright (c) 2019 C. Faloutsos 25

25

CarnegieMellon

Problem: e-commerce ratings

fraud

- User likes product
- User dislikes product
- User likes seller
- User dislikes seller
- Seller sells product
- Friends

Users
(honest, fraud)

Products
(high-quality-safe, low-quality-safe, non-malware, malware)

Sellers
(non-malware, malware)

- Given a heterogeneous graph on users, products, sellers and positive/negative ratings with “seed labels”
- Find the top k most fraudulent users, products and sellers

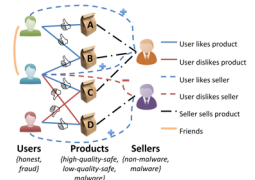
Dhivya Eswaran, Stephan Günnemann, Christos Faloutsos,
 “ZooBP: Belief Propagation for Heterogeneous Networks”,
 VLDB 2017

26

CarnegieMellon

Problem: e-commerce ratings fraud

DETAILS



Theorem 1 (ZooBP). *If $\mathbf{b}, \mathbf{e}, \mathbf{P}, \mathbf{Q}$ are constructed as described above, the linear equation system approximating the final node beliefs given by BP is:*

$$\mathbf{b} = \mathbf{e} + (\mathbf{P} - \mathbf{Q})\mathbf{b} \quad (\text{ZooBP}) \quad (10)$$

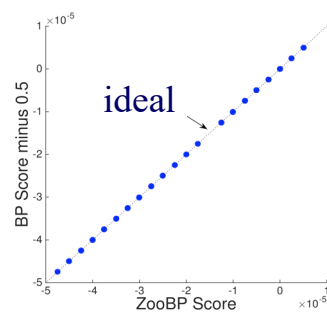
Dhivya Eswaran, Stephan Günnemann, Christos Faloutsos, “ZooBP: Belief Propagation for Heterogeneous Networks”, VLDB 2017

27

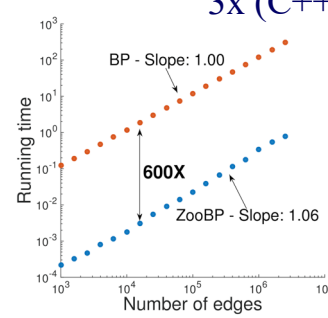
CarnegieMellon

ZooBP: features

Fast; convergence guarantees.



Near-perfect accuracy




linear in graph size

600x (matlab)
3x (C++)

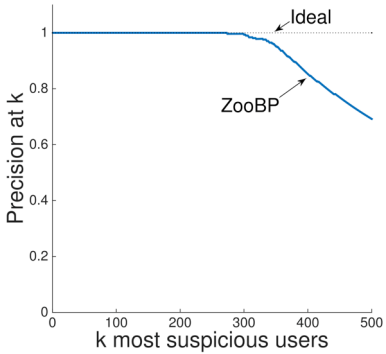
Dhivya Eswaran, Stephan Günnemann, Christos Faloutsos, “ZooBP: Belief Propagation for Heterogeneous Networks”, VLDB 2017

28

CarnegieMellon



ZooBP in the real world



k most suspicious users	Precision at k (Ideal)	Precision at k (ZooBP)
0	1.0	1.0
100	1.0	1.0
200	1.0	1.0
300	1.0	1.0
400	1.0	0.8
500	1.0	0.7

- Near 100% precision on top 300 users (Flipkart)
- Flagged users: suspicious
 - 400 ratings in 1 sec
 - 5000 good ratings and no bad ratings

Dhivya Eswaran, Stephan Günnemann, Christos Faloutsos, “ZooBP: Belief Propagation for Heterogeneous Networks”, VLDB 2017

29

CarnegieMellon

ZooBP: code etc

<http://www.cs.cmu.edu/~deswaran/code/zoobp.zip>



Dhivya Eswaran, Stephan Günnemann, Christos Faloutsos, “ZooBP: Belief Propagation for Heterogeneous Networks”, VLDB 2017

30

CarnegieMellon

Other ‘success stories’?

- Accounting fraud
- Malware detection

pwc

symantec

15-826 Copyright (c) 2019 C. Faloutsos 31

31

CarnegieMellon

Network Effect Tools: SNARE

- Some accounts are sort-of-suspicious – how to combine weak signals?

Before

pwc

Mary McGlohon, Stephen Bay, Markus G. Anderle, David M. Steier, Christos Faloutsos: *SNARE: a link analytic system for graph labeling and risk detection*. KDD 2009: 1265-1274


32

CarnegieMellon


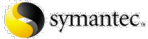
PATENT PENDING

Polonium: Tera-Scale Graph Mining and Inference for Malware Detection


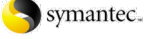
SDM 2011, Mesa, Arizona




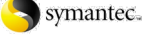
Polo Chau
Machine Learning Dept


Carey Nachenberg
Vice President & Fellow

Jeffrey Wilhelm
Principal Software Engineer

Adam Wright
Software Engineer




Prof. Christos Faloutsos
Computer Science Dept

33

CarnegieMellon

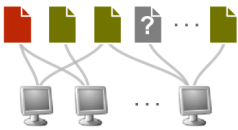
Polonium: The Data



60+ terabytes of data *anonymously* contributed by participants of worldwide *Norton Community Watch* program

50+ million machines

900+ million executable files



Constructed a machine-file bipartite graph (0.2 TB+)


1 billion nodes (machines and files)

37 billion edges

15-826
Copyright (c) 2019 C. Faloutsos
34


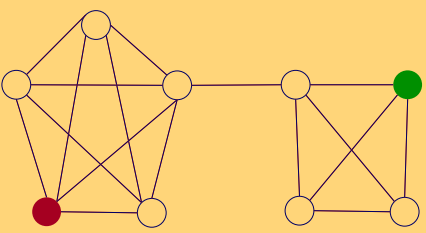
34

CarnegieMellon

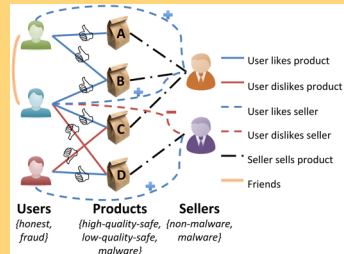


Short answer:

- What color, for the rest?
- A: Belief Propagation ('zooBP')

www.cs.cmu.edu/~deswaran/code/zoobp.zip



15-826 Copyright (c) 2019 C. Faloutsos 35