



Privacy Aspects of Social Graphs

Joseph Bonneau

Stanford Security Seminar, July 14 2009

University of Cambridge Computer Laboratory

Social Context And The Web



UNIVERSITY OF 800 YEARS CAMBRIDGE 1209~2009

Everything's Better With Friends...

- "Hyper-presence" of friends
- "networked public spaces"
- All web activity will have social context



Mike Barash wife just made pancakes and toast...not a bad way to start the day. and it appears we have power again. which is nice.

about an hour ago · Comment · Like



B[∞] 3 hours ago · Comment · Like

Griffin Barash is day 2 ... Options baby!



Ryan van Weezel at 2:27pm July 7 good luck... better have a red bull at lunch!

Adam Drewry at 4:36pm July 7 that sounds like a dream of a day

Write a comment...



Tyler Redlitz is celebrating his bday with beautiful weather in NYC!

5 hours ago · Comment · Like

📫 The Wu likes this.

Write a comment...



Facebook Is Becoming A Second Internet...

Function Page Markup **DB** Queries Email Forums Instant Messages **News Streams** Authentication Photo Sharing Video Sharing Blogging Microblogging **Micropayment Event Planning Classified Ads**

Internet version HTML, JavaScript SQL SMTP Usenet, etc. XMPP RSS OpenID Flickr, etc. YouTube, etc. Blogger, etc. Twitter, etc. Peppercoin, etc. E-Vite craigslist

Facebook version FBML FBQL **FB** Mail **FB** Groups FB Chat **FB** Stream **FB** Connect **FB** Photos FB Video **FB** Notes **FB** Status Updates **FB** Points **FB** Events FB Marketplace

Parallel Trend: The Internet is Becoming Social

"Given sufficient funding, all web sites expand in functionality until users can add each other as friends"



"Traditional" Social Network Analysis

- Performed by sociologists, anthropologists, etc. since the 70's
- Use data carefully collected through interviews & observation
 - Typically < 100 nodes
 - Complete knowledge
 - Links have consistent meaning
- All of these assumptions fail badly for online social network data



Traditional Graph Theory

- Nice Proofs
- Tons of definitions
- Ignored topics:
 - Large graphs
 - Sampling
 - Uncertainty



The Königsberg Bridges.

HAMILTON CYCLE ON DE BRUIJN GRAPH



Models Of Complex Networks From Math & Physics

Many nice models

- Erdos-Renyi
- Watts-Strogatz
- Barabasi-Albert
- Social Networks properties:
- Power-law
- Small-world
- High clustering coefficient



Real social graphs are complicated!



CAMBRIDGE 1209~2009

When In Doubt, Compute!

We do know many graph algorithms:

- Find important nodes
- Identify communities
- Train classifiers
- Identify anomalous connections

Major Privacy Implications!



• What can we infer purely from link structure?



 What can we infer purely from link structure? Joonwoong Kim A surprising amount! Andrew Lewis Popularity • Steven J. Murdoch Jonathan Anderson Centrality ٠ Luke Church Introvert vs. Extrovert ۲ Robert Watson Leadership potential ٠ Shailendra Fuloria Richard Clayton George Danezis

Markus Kuhn

• If we know nothing about a node but it's neighbours, what can we infer?

• If we know nothing about a node but it's neighbours, what can we infer?

A lot!

- Gender
- Political Beliefs
- Location
- Breed?

• Can we anonymise graphs?



- Can we anonymise graphs?
 - Not easily...
- Seminal result by Backstrom et al.: Attack of attack needs just 7 nodes
- Can do even better given user's complete neighborhood
- Also results for correlating users across networks
- Developing line of research...

• What can we infer if we "compromise" a fraction of nodes?



• What can we infer if we "compromise" a fraction of nodes?

A lot...

- Common theme: small groups of nodes can see the rest
 - Danezis et al.
 - Nagaraja
 - Korolova et al.
 - Bonneau et al.

- Can we defend against crawling in a sound way?
 - Work in progress!



• What if we get a subset of neighbours for all nodes?



• What if we get a subset of *k* neighbours for all nodes?

Emerging question for many social graphs

- Facebook and online SNS
- Mobile SNS



A Quietly Introduced Feature...



Facebook @ 2009 English (UK) \$

Log in About Advertising Developers Jobs Terms . Find Friends Privacy Help

Public Search Listings, Sep 2007

UNIVERSITY OF 800 YEARS CAMBRIDGE 1209~2009

Public Search Listings

- Unprotected against crawling
- Indexed by search engines
- Opt out—but most users don't know it exists!

Utility

Sign Up Sign up for	Facebook 1	o connect	t with Joe I	Bonneau.		Jeenneaale	, ginan com		
With the Joe Bonneau you were looking for? Search more	Joe Bor Add Joe Bo Here are so Dan Bragdon	meau as Frome of Joe B	riend Send J onneau's frie Corey Erickson	oe Bonneau ends: Jillian Day	a Message Vi	ew Joe Bonr	Bump Heldman	samantha Ricker	Joe Bonneau is on Facebook. Sign up for Facebook to connect with Joe Bonneau Sign Up It's free and anyone can join. Already a Member? Log to contact Joe Bonneau.
acebook © 2009 English (UK) \$					Log in a	About Adve	rtising Deve	elopers Jobs	Terms = Find Friends Priv







Promotion via Network Effects

UNIVERSITY OF 800 YEARS CAMBRIDGE 1209~2009



"Your name, network names, and profile picture thumbnail will be available in search results across the Facebook network and those limited pieces of information may be made available to third party search engines. This is primarily so your friends can find you and send a friend request."

-Facebook Privacy Policy



Legal Status



Much More Info Now Included...

UNIVERSITY OF 800 YEARS CAMBRIDGE 1209~2009

Legal Status

🔱 Advocates of Communism

Global

Basic Info

Type: Description:

Common Interest - Politics

The working class has nothing to lose but their chains. They have the world to win.

We have seen above that the first step in the revolution by the working class is to raise the proletariat to the position of ruling class to win the battle of democracy.

The proletariat will use its political supremacy to wrest, by degree, all capital from the bourgeoisie, to centralize all instruments of production in the hands of the state, i.e., of the proletariat organized as the ruling dass; and to increase the total productive forces as rapidly as possible.

Of course, in the beginning, this cannot be effected except by means of despotic inroads on the rights of property, and on the conditions of bourgeois production; by means of measures, therefore, which appear economically insufficient and untenable, but which, in the course of the movement, outstrip themselves, necessitate further inroads upon the old social order, and are unavoidable as a means of entirely revolutionizing the mode of production.

These measures will, of course, be different in different countries.

Nevertheless, in most advanced countries, the following will be pretty generally applicable.

1. Abolition of property in land and application of all rents of land to public purposes.

2. A heavy progressive or graduated income tax.

Members

Displaying 8 of 3,513 members





Group Type

This is an open group. Anyone can join and invite others to join.

Officers

Sim Party Philosopher

Nils Questioner of Party Authority

Aaron Official Representative of Anti-Revisionist Socialism

Aziz Official Representative of U.C.Y

William Official Representative of Moderate

Trotskyist Party Pmk Official Representative of Communist

Revolution Party Will

Official Representative of Utopia Party

Adem Vice Representative of the Downtrodden

Public Group Pages Recently Added

Obvious Attack

- Initially returned new friend set on refresh
- Can find all *n* friends in $O(n \cdot \log n)$ queries
 - The Coupon Collector's Problem
 - For 100 Friends, need 65 page refreshes
- As of Jan 2009, friends fixed per IP address

Fun with Tor

Germany

Australia

UNIVERSITY OF CAMBRIDGE

USA

UK



David Cottingham



Shoshana Freisinger



Eirik

Lauren Duffey



Emma

Conor Loftus-S





Raj



Stella

Nordhagen

Srilakshmi





Carl



Katie Gunderso n



David J

Hornsby

Sarita

Kristina

Sylvester





Ankit Garg





Brian Brown

Justin

Gary Champagn

Jillian





Cameron Lochte

Ben Skolnik



Melanie Kannokad

Shoshana

Freisinger

800 YEARS I 2 0 9 ~ 2 0 0 9

а



Federico

Baradello

Freisinger

Russ Heddlest

on

Lauren

Duffey



Luke

Church



Adrian Boscolo-





Hightower







Attack Scenario

- Spider all public listings
 - Our experiments crawled 250 k users daily
 - Implies ~800 CPU-days to recover all users

Abstraction

- Take a graph $G = \langle V, E \rangle$
- Randomly select k out-edges from each node

- Result is a sampled graph $G_k = \langle V, E_k \rangle$
- Try to approximate $f(G) \approx f_{approx}(G_k)$

Approximable Functions

- Node Degree
- Dominating Set
- Betweenness Centrality
- Path Length
- Community Structure

Experimental Data

- Crawled networks for Stanford, Harvard universities
- Representative sub-networks

	# Users	Mean d	Median d
Stanford	15043	125	90
Harvard	18273	116	76

Back To Our Abstraction

- Take a graph $G = \langle V, E \rangle$
- Randomly select k out-edges from each node

- Result is a sampled graph $G_k = \langle V, E_k \rangle$
- Try to approximate $f(G) \approx f_{approx}(G_k)$

Estimating Degrees

- Convert sampled graph into a directed graph
 - Edges originate at the node where they were seen
- Learn exact degree for nodes with degree < *k*
 - Less than k out-edges
- Get random sample for nodes with degree $\geq k$
 - Many have more than k in-edges


Average Degree: 3.5



Sampled with *k*=2





Degree known exactly for one node





Naïve approach: Multiply in-degree by average degree / k



Raise estimates which are less than *k*



Nodes with high-degree neighbors underestimated



Iteratively scale by current estimate / k in each step



After 1 iteration





Normalise to estimated total degree



Convergence after n > 10 iterations

- Converges fast, typically after 10 iterations
- Absolute error is high—38% average
 - Reduced to 23% for nodes with $d \ge 50$
- Still accurately can pick high degree nodes



Aggregate of x highest-degree nodes



Comparison of sampling parameters



• Set of Nodes $D \subseteq V$ such that

 $D \cup \text{Neighbours}(D) = V$

- Set allows viewing the entire network
- Also useful for marketing, trend-setting



Trivial Algorithm: Select High-Degree Nodes in Order



In fact, finding minimal dominating set is NP-complete



Greedy Algorithm: select for maximal coverage



Greedy Algorithm: select for maximal coverage



Shown to perform adequately in practice

Works Well on Sampled Graph



Insensitive to Sampling Parameter!



UNIVERSITY OF SCOYEARC CAMBRIDGE 1209-2009



- A measure of a node's importance
- Betweenness centrality:

$$C_B(v) = \sum_{s \neq v \neq t \in V} \frac{\sigma_{st}(v)}{\sigma_{st}}$$

• Measures the shortest paths in the graph that a particular vertex is part of

Centrality



- Goal: Find highly-connected sub-groups
- Measure success by high *modularity*:

$$Q = \frac{1}{2m} \sum_{v,w} \left[A_{vw} - \frac{d(v)d(w)}{2m} \right]$$

- Ratio of intra-community edges to random
- Normalised to be between -1 and 1



•Clausen et. al 2004 – find maximal modularity in $O(n \lg^2 n)$

•Track marginal modularity, update neighbours on each merge



























Conclusions

• *k*-sampling of each edge gives away a lot

Conclusions

• *k*-sampling of each edge gives away a lot

Can we fix it?



Regular subgraph extraction



Can we find a 2-regular subgraph?

Regular subgraph extraction



Step 1: Remove edges, weight by smallest attached node


Step 1: Remove edges, weight by smallest attached node



Step 1: Remove edges, weight by smallest attached node



Step 1: Remove edges, weight by smallest attached node



Step 1: Remove edges, weight by smallest attached node



Step 2: Remove further edges to force all degrees $\leq k$





Step 3: Randomly add edges between pairs of edges below k





Step 3: Randomly add edges between pairs of edges below k





(note: producing a cycle is atypical!)

How well have we done?

- Recall original goal of showing *k*-sample
 - Promotion, identification

- Two measures:
 - Precision: Percentage of edges shown which are real
 - *Recall*: Percentage of real edges which are shown (normalise recall to showing a max of *k* per node)

How well have we done?

- Recall original goal of showing *k*-sample
 - Promotion, identification

- Two measures:
 - Precision: Percentage of edges shown which are real
 - *Recall*: Percentage of real edges which are shown (normalise recall to showing a max of *k* per node)

	Original	Step 1	Step 2	Step 3
Precision	1	1	1	0.90
Recall	1	1	0.99	0.99





Drawbacks

- Requires complete graph knowledge
- Graph frequently changes!



Drawbacks

- Requires complete graph knowledge
- Graph frequently changes!

Alternative: Random Sampling

- Weight selection towards low-degree neighbours
- Computable locally, incrementally
- (much weaker...)

Random Sampling



Caveats

- Can gain some protection against degree estimation
 - With a lot of work
- Doesn't prevent inference of dominating sets, centrality!



Conclusions

- Availability of social graphs raises serious privacy concern
 - The blueprint of our society...
- Very fragile to many attacks
- Right now, we're choosing utility over privacy

Thank You!

jcb82@cl.cam.ac.uk

