Trace Complexity

Flavio Chierichetti Sapienza University of Rome

Online social networks



The total number of active user accounts exceeds two billions!



Advertisement



Note: includes advertising that appears on desktop and laptop computers as well as mobile phones and tablets, and includes all the various formats of advertising on those platforms; excludes SMS, MMS and P2P messaging-based advertising; numbers may not add up to 100% due to rounding Source: eMarketer, Sep 2012

Online Social Networks



Online Social Networks



Models of human behavior



Models Improvements

Models Validations

Online Social Networks



Models of human behavior



Empirical Analysis

- This sort of analysis is becoming more and more important in social networks research
- Exact user data is very hard to obtain for privacy related reasons

Empirical Analysis

- This sort of analysis is becoming more and more important in social networks research
- Exact user data is very hard to obtain for privacy related reasons
- Empirical researchers have a hard time testing their hypotheses / validating their models

"Data Traces"

Researchers try to **guess** the *missing data*, using *existing observations* and a reasonable *theoretical model*

"Data Traces"

Researchers try to **guess** the *missing data*, using *existing observations* and a reasonable *theoretical model*

Theoretical researchers can help by trying to understand whether **sound guessing algorithms** exist

Online Guessing Tasks

• Which of these people are friends?



Online Guessing Tasks

- Which of these people are friends?
- How many people were infected by this meme?



Online Guessing Tasks

- We see little pieces of some online social process.
- We would like to make some **sensible** guesses on the process as a whole.

Network Reconstruction

Joint with

Bruno Abrahao, Robert Kleinberg, Alessandro Panconesi

Network Reconstruction

- Incomplete Traces of Information Flow.
- Wish to infer the hidden network.

Network Reconstruction

- Rich area of study, pioneered by [Adar, Adamic, '05] in the context of social networks.
- MLE-based approaches:
 - [Gomez-Rodriguez, Leskovec, Krause '10],
 - [Gomez-Rodriguez, Balduzzi, Scholkopf '11],
 - [Myers, Leskovec '11],
 - [Du et al. '12]
- Information theoretic approaches
 - [Netrapalli, Sanghavi '12],
 - [Grippon, Rabbat '13]

"Inferring networks of diffusion and influence" Gomez-Rodriguez, Leskovec, Krause [KDD'10]

 Gomez-Rodriguez et al studied a large collection of blogs and memes to guess the *blogger network* that allowed *memes* to *spread*.



Abraham Lincoln invented Facebook MAY 08 2012

Intrigued by a possible connection between PT Barnum and Abe Lincoln, Nate St. Pierre travelled to the Lincoln Museum in Springfield, IL. Once there, he stumbled upon something called The Springfield Gazette, <u>a personal newspaper made by</u> Lincoln that is eerily similar to Facebook. "Inferring networks of diffusion and influence" Gomez-Rodriguez, Leskovec, Krause [KDD'10]

- Gomez-Rodriguez et al studied a large collection of blogs and memes to guess the *blogger network* that allowed *memes* to *spread*.
- They proposed a random meme-diffusion model, and used it for making the guess.

Alice's Blog

Bob's Blog - Live from Lewisville



The World According to Dave

Alice's Blog

Bob's Blog - Live from Lewisville

Mar 4, 2014, 8:25am Mrs. Hudson's cake shop will reopen!



The World According to Dave

Alice's Blog

Mar 4, 2014, 9:00am Rejoice! Mrs. Hudson cake shop is reopening

Bob's Blog - Live from Lewisville

Mar 4, 2014, 8:25am Mrs. Hudson's cake shop will reopen!



The World According to Dave

Alice's Blog

Mar 4, 2014, 9:00am Rejoice! Mrs. Hudson cake shop is reopening

Bob's Blog - Live from Lewisville

Mar 4, 2014, 8:25am Mrs. Hudson's cake shop will reopen!



charlie's blog

Mar 4, 2014, 10:00am OMG! Hudson's is back in business!

The World According to Dave

Alice's Blog	Bob's Blog - Live from Lewisville
<i>Mar 4, 2014, 9:00am</i>	<i>Mar 4, 2014, 8:25am</i>
Rejoice! Mrs. Hudson cake shop is reopening	Mrs. Hudson's cake shop will reopen!
C.	The World According to Dave
charlie's blog	My unsolicited ramblings, rants, musings & bloviations (totally free and worth every penny)
Mar 4, 2014, 10:00am	Mar 4, 2014, 10:30am
OMG! Hudson's is back in business!	Hudson's is about to reopen! Slurp!

Alice's Blog	Bob's Blog - Live from Lewisville
<i>Mar 4, 2014, 9:00am</i>	<i>Mar 4, 2014, 8:25am</i>
Rejoice! Mrs. Hudson cake shop is reopening	Mrs. Hudson's cake shop will reopen!
C.	The World According to Dave
charlíe's blog	My unsolicited ramblings, rants, musings & bloviations (totally free and worth every penny)
<i>Mar 4, 2014, 10:00am</i>	<i>Mar 4, 2014, 10:30am</i>
OMG! Hudson's is back in business!	Hudson's is about to reopen! Slurp!

Can we infer which blogger follows which blog?





















D 35' A 60' B 30' C 30' F 15' E



D 35' A 60' B 30' C 30' F 15' E B



B 10' E




B 10' E 20' F 15' A 15' D 30' C

Inference

Given a set of traces



Inference

Given a set of traces



we would like to infer the unknown blogger graph *G*

Gomez-Rodriguez, Leskovec, Krause [KDD'10]



• The unique source is chosen uniformly at random

Gomez-Rodriguez, Leskovec, Krause [KDD'10]



• The unique source is chosen uniformly at random



- The unique source is chosen uniformly at random
- The time to traverse an edge is an iid sample of $\mathrm{Exp}(\lambda)$



- The unique source is chosen uniformly at random
- The time to traverse an edge is an iid sample of $\mathrm{Exp}(\lambda)$



- The unique source is chosen uniformly at random
- The time to traverse an edge is an iid sample of $\mathrm{Exp}(\lambda)$
- The trace follows the shortest path tree from the source



- The unique source is chosen uniformly at random
- The time to traverse an edge is an iid sample of $\mathrm{Exp}(\lambda)$
- The trace follows the shortest path tree from the source



- The unique source is chosen uniformly at random
- The time to traverse an edge is an iid sample of $\mathrm{Exp}(\lambda)$
- The trace follows the shortest path tree from the source



- The unique source is chosen uniformly at random
- The time to traverse an edge is an iid sample of $\mathrm{Exp}(\lambda)$
- The trace follows the shortest path tree from the source



- The unique source is chosen uniformly at random
- The time to traverse an edge is an iid sample of $\mathrm{Exp}(\lambda)$
- The trace follows the shortest path tree from the source



- The unique source is chosen uniformly at random
- The time to traverse an edge is an iid sample of $\mathrm{Exp}(\lambda)$
- The trace follows the shortest path tree from the source

Gomez-Rodriguez, Leskovec, Krause [KDD'10]



The unique source is chosen uniformly at random

- The time to traverse an edge is an iid sample of $\mathrm{Exp}(\lambda)$
- The trace follows the shortest path tree from the source

Gomez-Rodriguez, Leskovec, Krause [KDD'10]



The unique source is chosen uniformly at random

- The time to traverse an edge is an iid sample of $\mathrm{Exp}(\lambda)$
- The trace follows the shortest path tree from the source



- The unique source is chosen uniformly at random
- The time to traverse an edge is an iid sample of $\mathrm{Exp}(\lambda)$
- The trace follows the shortest path tree from the source

Can we make a sensible guess of the network?

Can we make a sensible guess of the network?

How many traces would we need?

Abrahao, Chierichetti, Kleinberg, Panconesi [KDD'13]

- We show that $O\left(n\Delta\log n\right)$ traces are sufficient for reconstruction.
- We show that, in some cases, $\Omega\left(\frac{n\Delta}{\log^2 n}\right)$ traces are necessary.

n is the number of nodes of the unknown graph, Δ is its maximum degree.

- We show that $O\left(n\Delta\log n\right)$ traces are sufficient for reconstruction.
- We show that, in some cases, $\Omega\left(\frac{n\Delta}{\log^2 n}\right)$ traces are necessary.
- We also show that $O(\operatorname{poly}(\Delta) \cdot \log n)$ traces are sufficient, that $O(\log n)$ traces are sufficient for trees...

Observation

If u and v are the first two nodes of a trace then the edge $\{u,v\}$ is in the unknown graph.



Observation

If u and v are the first two nodes of a trace then the edge $\{u,v\}$ is in the unknown graph.



Observation

If u and v are the first two nodes of a trace then the edge $\{u, v\}$ is in the unknown graph.

 We can then perfectly reconstruct if, for each unknown edge {u,v}, there exists a trace that begins with its two endpoints.



Observation

If u and v are the first two nodes of a trace then the edge $\{u, v\}$ is in the unknown graph.

 We can then perfectly reconstruct if, for each unknown edge {u,v}, there exists a trace that begins with its two endpoints.



Observation

If u and v are the first two nodes of a trace then the edge $\{u, v\}$ is in the unknown graph.

 We can then perfectly reconstruct if, for each unknown edge {u,v}, there exists a trace that begins with its two endpoints.

- The probability that a specific node is the first one in a trace is 1/n.
- The probability that the edge {u,v} is spanned by the first two nodes of the trace is then at least $\frac{1}{n \cdot \min(\deg(u), \deg(v))}$

- The probability that a specific node is the first one in a trace is 1/n.
- The probability that the edge {u,v} is spanned by the first two nodes of the trace is then at least



- The probability that a specific node is the first one in a trace is 1/n.
- The probability that the edge {u,v} is spanned by the first two nodes of the trace is then at least $\frac{1}{n \cdot \min(\deg(u), \deg(v))}$
- Any classical tail bound then proves that $O(n\Delta \log n)$ traces are enough to perfectly reconstruct the graph.

The first-edge algorithm

- Works for any (uniform) edge waiting time distribution.
- It is competitive with more complex machinelearning based algorithms (Gomez-Rodriguez, Leskovec, Krause [KDD'10]).

Performances



- We show that $O\left(n\Delta\log n\right)$ traces are sufficient for reconstruction.
- We show that, in some cases, $\Omega\left(\frac{n\Delta}{\log^2 n}\right)$ traces are necessary.

- We show that $O(nD\log n)$ traces are sufficient for guessing the edges incident on a node of degree $\leq D$
- We show that, in some cases, $\Omega\left(\frac{n\Delta}{\log^2 n}\right)$ traces are necessary.

- We show that $O(nD\log n)$ traces are sufficient for guessing the edges incident on a node of degree $\leq D$
- We show that, in some cases, $\Omega\left(\frac{n\Delta}{\log^2 n}\right)$ traces are necessary.

- We show that $O(nD\log n)$ traces are sufficient for guessing the edges incident on a node of degree $\leq D$
- We show that, in some cases, $\Omega\left(\frac{n\Delta}{\log^2 n}\right)$ traces are necessary.

Why ~ 40K traces?

Most edges are incident on nodes having small degree


Why ~ 40K traces?

Most edges are incident on nodes having small degree



Back-of-envelope calculation

Why ~ 40K traces?

Most edges are incident on nodes having small degree



Why ~ 40K traces?

Most edges are incident on nodes having small degree



Too many traces?

- In some cases such a large number of traces is not available.
- In those cases, one can still try to reconstruct some (less precise) missing information about the network.

E-mail Activisim

Joint with Jon Kleinberg, David Liben-Nowell

Internet Activism

- Very important phenomenon
- Incomplete Traces
- Chain Letter Petitions: how to estimate the reach?

NPR Chain Letter

PBS, NPR (National Public Radio), and the arts are facing major cutbacks in funding. In spite of the efforts of each station to reduce spending costs and streamline their services, the government officials believe that the funding currently going to these programs is too large a portion of funding for something which is seen as "unworthwhile."

[...]

When this issue comes up in 1996, the funding will be determined for fiscal years 1996-1998.

The only way that our representatives can be aware of the base of support or PBS and funding for these types of programs is by making our voices heard.

Please add your name to this list if you believe in what we stand for. This list will be forwarded to the President of the United States, the Vice President of the United States, the House of Representatives and Congress.

If you happen to be the 50th, 100th, 150th, etc. signer of this petition, please forward to: kubi7975@blue.univnorthco.edu . This way we can keep track of the lists and organize them. Forward this to everyone you know, and help us to keep these programs alive.

Thank you.

1. Elizabeth Weinert, student, University of Northern Colorado, Greeley, Colorado. 2. Robert M. Penn; San Francisco, CA 3. Gregory S. Williamson, San Francisco, CA 4. Daniel C. Knightly, Austin, TX 5. Andrew H. Knightly, Los Angeles, CA 6. Aaron C. Yeater, Somerville, MA 7. Tobie M. Cornejo, Washington, DC 8. John T. Mason, Dalton, MA 9. Eric W. Fish, Williamstown, MA 10. Courtney E. Estill, Hamilton College, NY 11. Vanessa Moore, Northfield, MN 12. Lynne Raschke, Haverford College, PA (originally Minnesota) 13. Deborah Bielak, Haverford, PA 14. Morgan Lloyd, Haverford, PA 19041 15. Galen Lloyd, Goucher College, MD 16. Brian Eastwood, University of Vermont, VT 17. Elif Batuman, Harvard University, MA 18. Kohar Jones, Yale University, CT 19. Claudia Brittenham, Yale University, CT 20. Alexandra Block, Yale University, CT 21. Susanna Chu, Yale University, CT 22. Michelle Chen, Harvard University, MA 23. Jessica Hammer, Harvard University, MA 24. Ann Pettigrew, Haverford College, PA 25. Kirstin Knox, Swarthmore College, PA 26. Jason Adler, Swarthmore College, PA 27. Daniel Gottlieb, Swarthmore College (but truly from Lawrence, KS) 28. Josh Feltman, Tufts University, MA 29. Louise Forrest, Massachusetts Institute of Technology, MA 30. HongSup Park, Massachusetts Institute of Technology, MA (originally from Portage, Wisconsin) 31. Ana Sandoval, Massachusetts Institute of Technology [...]

NPR Chain Letter

PBS, NPR (National Public Radio), and the arts are facing major cutbacks in funding. In spite of the efforts of each station to reduce spending costs and streamline their services, the government officials believe that the funding currently going to these programs is too large a portion of funding for something which is seen as "unworthwhile."

[...]

When this issue comes up in 1996, the funding will be determined for fiscal years 1996-1998.

The only way that our representatives can be aware of the base of support or PBS and funding for these types of programs is by making our voices heard.

Please add your name to this list if you believe

<u>in what we stand for</u>. This list will be forwarded to the President of the United States, the Vice President of the United States, the House of Representatives and Congress.

If you happen to be the 50th, 100th, 150th, etc. signer of this petition, please forward to: kubi7975@blue.univnorthco.edu . This way we can keep track of the lists and organize them. <u>Forward</u> <u>this to everyone you know, and help us to keep</u> these programs alive.

Thank you.

1. Elizabeth Weinert, student, University of Northern Colorado, Greeley, Colorado. 2. Robert M. Penn; San Francisco, CA 3. Gregory S. Williamson, San Francisco, CA 4. Daniel C. Knightly, Austin, TX 5. Andrew H. Knightly, Los Angeles, CA 6. Aaron C. Yeater, Somerville, MA 7. Tobie M. Cornejo, Washington, DC 8. John T. Mason, Dalton, MA 9. Eric W. Fish, Williamstown, MA 10. Courtney E. Estill, Hamilton College, NY 11. Vanessa Moore, Northfield, MN 12. Lynne Raschke, Haverford College, PA (originally Minnesota) 13. Deborah Bielak, Haverford, PA 14. Morgan Lloyd, Haverford, PA 19041 15. Galen Lloyd, Goucher College, MD 16. Brian Eastwood, University of Vermont, VT 17. Elif Batuman, Harvard University, MA 18. Kohar Jones, Yale University, CT 19. Claudia Brittenham, Yale University, CT 20. Alexandra Block, Yale University, CT 21. Susanna Chu, Yale University, CT 22. Michelle Chen, Harvard University, MA 23. Jessica Hammer, Harvard University, MA 24. Ann Pettigrew, Haverford College, PA 25. Kirstin Knox, Swarthmore College, PA 26. Jason Adler, Swarthmore College, PA 27. Daniel Gottlieb, Swarthmore College (but truly from Lawrence, KS) 28. Josh Feltman, Tufts University, MA 29. Louise Forrest, Massachusetts Institute of Technology, MA 30. HongSup Park, Massachusetts Institute of Technology, MA (originally from Portage, Wisconsin) 31. Ana Sandoval, Massachusetts Institute of Technology [...]





























George and Hilary, by exposing their emails, revealed a subtree of the Chain Letter tree.

Real-World Chain Letters' Tree

- [Liben-Nowell, Kleinberg, PNAS'08], mined
 - web-accessible mailing-lists,
 - blog posts.
- They obtained some "exposed" nodes of two Chain Letters' trees, and
- they produced two "revealed" trees.

NPR revealed tree

Liben-Nowell, Kleinberg, PNAS'08



NPR revealed tree Liben-Nowell, Kleinberg, PNAS'08

		g	800	
		88	5666	è.
	ø	~\$ \$, de
	độ	\$ 8		\$ \$
	- \$ \$	\$ '	• <u>0 0 0</u>	8
	- \$\$	\$	\$\$\$; •
	33	\$	\$	5
ç	899 89	\$	ခွဲခွဲ	•
ę ę	\$\$	٠	800	è.
Ş	\$\$		\$\$	\$
\$	\$\$		\$\$	\$
~ &	- \$\$		<u>\$</u> \$	\$
\$\$	- \$\$		\$\$	\$
\$ \$	÷\$	2	500 500	99
666 666	- \$\$	2		<u></u>
\$\$\$	33	- 2	300	\$\$
\$\$\$	866 866		333	\$\$
\$\$\$	<u> </u>	ž	333	\$\$
\$\$\$	900 600	୍ବ		<u> </u>
\$\$\$	ခွံခံခွဲ	- \$3		\$
\$\$\$	\$ \$	- 33		\$
\$ \$ \$	\$ \$	୍ଚିଶ୍	\$\$\$\$	\$
\$ <u>\$</u> \$	\$ \$	333		\$
ဓဓဓ	နို နိုင်	200		Ş.
ဓဓဓ	မှ မို	3999	နှင့်ခွင့်	•
ဓဓဓ	မှ မို	200	ဓဓဓဓ	
ဓိဓိဓ	နှိန့်	₹ •7		•
ဒိုမိုမို	နိုန်	33		•
ဓဓဓ	မှ မှ	4	999	•
ဓဓဓ ဓဓဓ	9 • (4	999	
ဓဓဓ ဓဓဓ	ş ç	4	မိုမို	
ဓဓဓ	ş	5 5	•	
333	8 8	33	3	

NPR revealed tree Liben-Nowell, Kleinberg, PNAS'08

00000000

			29	የድዮያ	2
			23	3225	šo-
		8	φ¢	şęęć	φ.
	5	22 1	23	2229	22
	2	53 3	3	1444	56
		55 (φ.	999	5ō-
	5	<u>è</u>	<u> </u>	- 999	•
	- 2	88 3	2	223	2
	2	53 1	Υ.	221	~
	ø	è	¢.	- ¢¢•	•
	2 2	22	9	22	à
	58	55	-	201	5
<	5 6	50		φφ	Ŷ.
5	2 2	29		<u> </u>	<u> </u>
2	53	32		22	8
3	६ २	58 -		22	ę.
99	2 4	20		<u> </u>	¢.
- 22	2 2	22		22	2
- 2,5	52	53		322	20
¢ ģ	\$ \$	è 🔶		¢¢¢	φφ
999	2 9	29		200	99
- 220	52	55	- 2	555	22
- 444	53	55		794	99
<u> </u>	2 99	24		<u> </u>	$\dot{\gamma}\dot{\gamma}$
- 222	2 X X	22	- 2	222	22
- 225	5 53	55	3	555	99
¢ ¢ ¢	è 🔆	2¢	ge	è è è	 • • •
222	2 23	22	23	5.22	2°
- 222	5 8	4	- 29	5666	5
¢¢¢	è 🔶	¢.	- \$9	<u>è</u> фф	φ.
222	22	2	22	2222	22
- 22	58	50	530	5225	5
	5 ¢	φş	5 ф	5994	φ.
222	22	22	239	2222	<u>9</u>
- 223	53	330	325	5555	32
- ġġa	5 ō	φφ	φφ¢	5999	•
999	2 2	999	20	• • • • •	2
220	5 2	220	32	666	5
- 444	5 8	- 44	ΫŶ.	999	
	2 9	<u> </u>	ŧ¢.	<u> </u>	•
223	22	22	2	222	2
- 223	53	22	2	225	\$
- 666	ģ ģ	φφ	٠.	Ģ Ģē	>
999	29	٠ģ	Ŷ	222	2
666	58	2	3	200	5
- 444	5 ŏ	- Ą	4	444	
	è è	- †	\$	ŧф	
222	22	2	2	2	
- 223	58	2	2	2	
1 1 1		I	I		





Iraq Chain Letter

Dear all:

The US Congress has just authorized the President of the US to go to war against Iraq. The UN is gathering signatures in an effort to avoid this tragic world event.

Please consider this an urgent request: UN Petition for Peace - Stand for Peace. Islam is not the Enemy. War is NOT the Answer.

Today we are at a point of imbalance in the world and are moving toward what may be the beginning of a THIRD WORLD WAR.

Please COPY (rather than Forward) this e-mail in a new message, sign at the end of the list, and send it to all the people whom you know.

If you receive this list with more than 500 names signed, please send a copy of the message to:

usa@un.int president@whitehouse.gov

Even if you decide not to sign, please consider forwarding the petition on instead of deleting it.

Suzanne Dathe, Grenoble, France
 Laurence COMPARAT, Grenoble, France
 Philippe MOTTE, Grenoble, France

4) Jok FERRAND, Mont St. Martin, France 5) Emmanuelle PIGNOL, St Martin d'Heres, FRANCE 6) Marie GAUTHIER, Grenoble, FRANCE 7) Laurent VESCALO, Grenoble, FRANCE 8) Mathieu MOY, St Egreve, FRANCE 9) Bernard BLANCHET, Mont St Martin, FRANCE 10) Tassadite FAVRIE, Grenoble, FRANCE 11) Loic GODARD, St Ismier, FRANCE 12) Benedicte PASCAL, Grenoble, FRANCE 13) Khedaidja BENATIA, Grenoble, FRANCE 14) Marie-Therese LLORET, Grenoble, FRANCE 15) Benoit THEAU, Poitiers, FRANCE 16) Bruno CONSTANTIN, Poitiers, FRANCE 17) Christian COGNARD, Poitiers, FRANCE 18) Robert GARDETTE, Paris, FRANCE 19) Claude CHEVILLARD, Montpellier, FRANCE 20) Gilles FREISS, Montpellier, FRANCE 21) Patrick AUGEREAU, Montpellier, FRANCE 22) Jean IMBER! T, Marseille, FRANCE 23) Jean-Claude MURAT, Toulouse, France 24) Anna BASSOLS, Barcelona, Catalonia 25) Mireia DUNACH, Barcelona, Catalonia 26) Michel VILLAZ, Grenoble, France 27) Pages Frederique, Dijon, France 28) Rodolphe FISCHMEISTER, Chatenay-Malabry, France 29) Francois BOUTEAU, Paris, France 30) Patrick PETER, Paris, France 31) Lorenza RADICI, Paris, France 32) Monika Siegenthaler, Bern, Switzerland 33) Mark Philp, Glasgow, Scotland 34) Tomas Andersson, Stockholm, Sweden 35) Jonas Eriksson, Stockholm, Sweden 36) Karin Eriksson, Stockholm, Sweden . . .

Iraq Chain Letter

Dear all:

The US Congress has just authorized the President of the US to go to war against Iraq. The UN is gathering signatures in an effort to avoid this tragic world event.

Please consider this an urgent request: UN Petition for Peace - Stand for Peace. Islam is not the Enemy. War is NOT the Answer.

Today we are at a point of imbalance in the world and are moving toward what may be the beginning of a THIRD WORLD WAR.

Please COPY (rather than Forward) this e-mail in a new message, sign at the end of the list, and send it to all the people whom you know.

If you receive this list with more than 500 names signed, please send a copy of the message to:

usa@un.int president@whitehouse.gov

Even if you decide not to sign, please consider forwarding the petition on instead of deleting it.

Suzanne Dathe, Grenoble, France
 Laurence COMPARAT, Grenoble, France
 Philippe MOTTE, Grenoble, France

4) Jok FERRAND, Mont St. Martin, France 5) Emmanuelle PIGNOL, St Martin d'Heres, FRANCE 6) Marie GAUTHIER, Grenoble, FRANCE 7) Laurent VESCALO, Grenoble, FRANCE 8) Mathieu MOY, St Egreve, FRANCE 9) Bernard BLANCHET, Mont St Martin, FRANCE 10) Tassadite FAVRIE, Grenoble, FRANCE 11) Loic GODARD, St Ismier, FRANCE 12) Benedicte PASCAL, Grenoble, FRANCE 13) Khedaidja BENATIA, Grenoble, FRANCE 14) Marie-Therese LLORET, Grenoble, FRANCE 15) Benoit THEAU, Poitiers, FRANCE 16) Bruno CONSTANTIN, Poitiers, FRANCE 17) Christian COGNARD, Poitiers, FRANCE 18) Robert GARDETTE, Paris, FRANCE 19) Claude CHEVILLARD, Montpellier, FRANCE 20) Gilles FREISS, Montpellier, FRANCE 21) Patrick AUGEREAU, Montpellier, FRANCE 22) Jean IMBER! T, Marseille, FRANCE 23) Jean-Claude MURAT, Toulouse, France 24) Anna BASSOLS, Barcelona, Catalonia 25) Mireia DUNACH, Barcelona, Catalonia 26) Michel VILLAZ, Grenoble, France 27) Pages Frederique, Dijon, France 28) Rodolphe FISCHMEISTER, Chatenay-Malabry, France 29) Francois BOUTEAU, Paris, France 30) Patrick PETER, Paris, France 31) Lorenza RADICI, Paris, France 32) Monika Siegenthaler, Bern, Switzerland 33) Mark Philp, Glasgow, Scotland 34) Tomas Andersson, Stockholm, Sweden 35) Jonas Eriksson, Stockholm, Sweden 36) Karin Eriksson, Stockholm, Sweden . . .



IRAQ revealed tree Liben-Nowell, Kleinberg, PNAS'08

18,119 nodes
17,079 nodes with one child (94%)

IRAQ revealed tree

Liben-Nowell, Kleinberg, PNAS'08

18,119 nodes
17,079 nodes with one child (94%)
620 exposed nodes
557 (exposed) leaves

IRAQ revealed tree

Liben-Nowell, Kleinberg, PNAS'08

18,119 nodes
17,079 nodes with one child (94%)
620 exposed nodes
557 (exposed) leaves

Why is this fraction so high?

IRAQ revealed tree

Liben-Nowell, Kleinberg, PNAS'08

18,119 nodes
17,079 nodes with one child (94%)
620 exposed nodes
557 (exposed) leaves

Why is this fraction so high?

What can we infer about the original, *unknown*, Chain Letter Tree?

Tree-Revealing Process Liben-Nowell, Kleinberg, PNAS'08



Tree-Revealing Process Liben-Nowell, Kleinberg, PNAS'08



Each node is exposed independently with prob. $\delta > 0$

Tree-Revealing Process Liben-Nowell, Kleinberg, PNAS'08



Each node is exposed independently with prob. $\delta > 0$


Each node is exposed independently with prob. $\delta > 0$



Each node is exposed independently with prob. $\delta > 0$



Ancestors of exposed nodes are revealed



Ancestors of exposed nodes are revealed

Previous Work

- Golub, Jackson, PNAS'10 perform simulations,
 - using branching process trees near the critical threshold as the Chain Letter Trees,
 - and exposing nodes as in Kleinberg, Liben-Nowell, PNAS'08.
- They observe that the *revealed* tree has a high fraction of nodes with only one child (and some other properties).

 Our 1st result, informally, states that the tree-revealing process, is enough to explain the high fraction of single-child nodes

 Our 1st result, informally, states that the tree-revealing process, is enough to explain the high fraction of single-child nodes assuming only a degree bound on the unknown chain letter tree.

We see a "revealed" tree...



We see a "revealed" tree...



...we would like to study the "unknown" tree!



We see a "revealed" tree...



...we would like to study the "unknown" tree!



Size? Width? Height? Degree Distribution? ...

We see a "revealed" tree...



...we would like to study the "unknown" tree!



Size? Width? Height? Degree Distribution? ...

 Our 2nd result, informally, states that (under reasonable assumptions) it is possible to estimate the size of the unknown chain letter tree with a small error, with high probability.

 Our 2nd result, *informally*, states that (under reasonable assumptions) it is possible to estimate the size of the unknown chain letter tree with a small error, with high probability.

Observe that we do not know the exposing probability δ

 Our 2nd result, informally, states that (under reasonable assumptions) it is possible to estimate the size of the unknown chain letter tree with a small error, with high probability.

We use this theorem to estimate that ~ 173k people that signed the IRAQ chain letter

This estimate is backed by a probability bound (on the probability space induced by the revealing process)

 Our 2nd result, *informally*, states that (under reasonable assumptions) it is possible to estimate the size of the unknown chain letter tree with a small error, with high probability.

We use this theorem to estimate that ~ 173k people that signed the IRAQ chain letter

The chain letter generated ~ 3.5M emails

Single-Child Fraction

- Nodes are exposed with probability $\delta > 0$
- We assume that the unknown tree's maximum degree is at most k



































in such a way that each subforest has $\simeq \delta^{-1}$ nodes and the median height in the subforest is $\Omega(\log_{k-1} \delta^{-1})$.



 $\Pr[\Omega(n \cdot \delta \cdot \log_{k-1} \delta^{-1}) \text{ nodes will be revealed}] = 1 - o(1)$

Single-Child Fraction

 $\Pr[\Omega(n \cdot \delta \cdot \log_{k-1} \delta^{-1}) \text{ nodes will be revealed}] = 1 - o(1)$

Single-Child Fraction

 $\Pr[\Omega(n \cdot \delta \cdot \log_{k-1} \delta^{-1}) \text{ nodes will be revealed}] = 1 - o(1)$

 $\Pr[\text{at most } 2 \cdot n \cdot \delta \text{ nodes will be exposed}] = 1 - o(1)$
$\Pr[\Omega(n \cdot \delta \cdot \log_{k-1} \delta^{-1}) \text{ nodes will be revealed}] = 1 - o(1)$

 $\Pr[\text{at most } 2 \cdot n \cdot \delta \text{ nodes will be exposed}] = 1 - o(1)$

Each leaf in the revealed tree is an exposed node.



 $\Pr[\Omega(n \cdot \delta \cdot \log_{k-1} \delta^{-1}) \text{ nodes will be revealed}] = 1 - o(1)$

 $\Pr[\text{at most } 2 \cdot n \cdot \delta \text{ nodes will be exposed}] = 1 - o(1)$

Each leaf in the revealed tree is an exposed node.

 $\Pr[\text{the revealed tree will have at most } 2 \cdot n \cdot \delta \text{ leaves}] = 1 - o(1)$

 $\Pr[\Omega(n \cdot \delta \cdot \log_{k-1} \delta^{-1}) \text{ nodes will be revealed}] = 1 - o(1)$

 $\Pr[\text{at most } 2 \cdot n \cdot \delta \text{ nodes will be exposed}] = 1 - o(1)$

Each leaf in the revealed tree is an exposed node.

 $\Pr[\text{the revealed tree will have at most } 2 \cdot n \cdot \delta \text{ leaves}] = 1 - o(1)$

In an arbitrary tree, the number of internal nodes with more than one child is upper-bounded by the number of leaves.

 $\Pr[\Omega(n \cdot \delta \cdot \log_{k-1} \delta^{-1}) \text{ nodes will be revealed}] = 1 - o(1)$

 $\Pr[\text{at most } 2 \cdot n \cdot \delta \text{ nodes will be exposed}] = 1 - o(1)$

Each leaf in the revealed tree is an exposed node.

 $\Pr[\text{the revealed tree will have at most } 2 \cdot n \cdot \delta \text{ leaves}] = 1 - o(1)$

In an arbitrary tree, the number of internal nodes with more than one child is upper-bounded by the number of leaves.

 $\Pr[\Omega(n \cdot \delta \cdot \log_{k-1} \delta^{-1}) \text{ nodes will be revealed}] = 1 - o(1)$

 $\Pr[\Omega(n \cdot \delta \cdot \log_{k-1} \delta^{-1}) \text{ nodes will be revealed}] = 1 - o(1)$



 $\Pr[\Omega(n \cdot \delta \cdot \log_{k-1} \delta^{-1}) \text{ nodes will be revealed}] = 1 - o(1)$

 $\Pr[\Omega(n \cdot \delta \cdot \log_{k-1} \delta^{-1}) \text{ nodes will be revealed}] = 1 - o(1)$

 $\Pr[\Omega(n \cdot \delta \cdot \log_{k-1} \delta^{-1}) \text{ nodes will be revealed}] = 1 - o(1)$



 $\Pr[\Omega(n \cdot \delta \cdot \log_{k-1} \delta^{-1}) \text{ nodes will be revealed}] = 1 - o(1)$

$$\mathsf{A} \ \frac{1}{\log_{k-1} \delta^{-1}} \ \text{fraction of the set.}} \gg n \cdot \delta \cdot \log_{k-1} \delta^{-1}$$

 $\Pr[\Omega(n \cdot \delta \cdot \log_{k-1} \delta^{-1}) \text{ nodes will be revealed}] = 1 - o(1)$

 $\Pr[\text{the revealed tree has} \le 4n\delta \text{ non-single-child nodes}] = 1 - o(1)$

Pr [the fraction of single-child nodes in the

revealed tree is
$$\geq 1 - O\left(\frac{1}{\log_{k-1}\delta^{-1}}\right) = 1 - o(1)$$

 $\Pr[\Omega(n \cdot \delta \cdot \log_{k-1} \delta^{-1}) \text{ nodes will be revealed}] = 1 - o(1)$

 $\Pr[\text{the revealed tree has} \le 4n\delta \text{ non-single-child nodes}] = 1 - o(1)$

Pr [the fraction of single-child nodes in the

revealed tree is
$$\geq 1 - O\left(\frac{1}{\log_{k-1}\delta^{-1}}\right) = 1 - o(1)$$

The high single-child fraction can be explained by assuming just a degree bound on the unknown tree

Number of Signers

Revealed Tree

Unknown Tree





How to guess the size of the unknown tree?





















Revealed Tree



Revealed Tree



















3. Take the ratio



What can go wrong?








Theorem

 The previous algorithm can guess the size with high probability if

$$n > \tilde{\Omega}\left(\max\left(\delta^{-2}, \delta^{-1} \cdot k \right) \right)$$

- k is the maximum number of children in the unknown tree,
- δ is the exposing probability.

Theorem

 The previous algorithm can guess the size with high probability if

$$n > \tilde{\Omega}\left(\max\left(\delta^{-2}, \delta^{-1} \cdot k \right) \right)$$

- k is the maximum number of children in the unknown tree,
- δ is the exposing probability.

$$k < \tilde{O}\left(\sqrt{n}\right)$$
 $\delta > \tilde{\Omega}\left(\sqrt{\frac{1}{n}}\right)$ satisfy the requirement

Theorem

 The previous algorithm can guess the size with high probability if

$$n > \tilde{\Omega}\left(\max\left(\delta^{-2}, \delta^{-1} \cdot k \right) \right)$$

- k is the maximum number of children in the unknown tree,
- δ is the exposing probability.
- No algorithm can do it if *n* is smaller.

IRAQ Tree Size

- We refined our asymptotic theorem for the IRAQ revealed tree (18k nodes)
- Assuming the tree-revealing model, we estimate that the number of signers of the IRAQ petition is within a factor of 2 of 173k with probability ≥ 95%

Conclusion

Reconstruction Problems

- Answers depend on:
 - the random model that drives the process;
 - the number of traces we have;
 - the random algorithm that reveals pieces of information by "cutting" them out of a model run.

Random Model

- It is important to check whether the random model is close to reality;
- without this step, the guesses might be "close" to the model's runs, but very far from reality.

This step is possibly the hardest one to do.

Random Model

- It is important to check whether the random model is close to reality;
- without this step, the guesses might be "close" to the model's runs, but very far from reality.

This step is possibly the hardest one to do. Checking whether the guesses are close to reality is sometimes

the only viable approach.

Number of Traces

- If we want to get a significant understanding of some partially-known process,
- we first and foremost need to verify whether the questions we are asking can be significantly answered by the amount of data we have.

This step can be carried out (theorems and/or simulations)

Trace Generation

- The number of "traces" needed depends strongly on which pieces of information of a model's run are revealed (as well as on the random model).
- In some cases, it might be possible to get more informative traces by a deeper mining of the data.

Complexity of Guessing

- If not enough traces are available, we should simplify the questions we ask about the unknown process.
- Otherwise, our guesses might be completely off and insignificant.

Complexity of Guessing

- How many traces do we need for other guessing tasks?
- Very, very, rich area of problems.

Thanks!