Mining Social and Semantic Network Data on the Web



Markus Schatten, PhD

University of Zagreb Faculty of Organization and Informatics

May 4, 2011

Introduction

Web 2.0, Semantic Web, Web 3.0

Network science

Data collection

Semistructured data XPath Regular expressions Ajax APIs

Data analysis

Network matrices Folksonomy model Logic programming

Case study - CROSBI



2 of 54

Outline



Introduction

Web 2.0, Semantic Web, Web 3.0

Network science

Data collection

Data analysis

Case study - CROSBI

Questions



- Why mining social and semantic network data?
- How to collect these data?
- What are the common obstacles and how to solve them?
- How to analyze these data?

Outline



Introduction

Web 2.0, Semantic Web, Web 3.0

Network science

Data collection

Data analysis

Case study - CROSBI

Web 2.0, Semantic Web, Web 3.0





Why should one bother with these buzzwords?



Millions of people,



Millions of people, on thousands of sites,



Millions of people, on thousands of sites, across hundreds of diverse cultures,



Millions of people, on thousands of sites, across hundreds of diverse cultures, leave "trails"



Millions of people, on thousands of sites, across hundreds of diverse cultures, leave "trails" about various aspects of their lives...



Millions of people, on thousands of sites, across hundreds of diverse cultures, leave "trails" about various aspects of their lives...

"Social systems are systems of communication"

(Niklas Luhman)



Millions of people, on thousands of sites, across hundreds of diverse cultures, leave "trails" about various aspects of their lives...

"Social systems are systems of communication"

(Niklas Luhman)

By analyzing the "trails" (the results of social systems structural coupling) we can analyze the social system itself:



Millions of people, on thousands of sites, across hundreds of diverse cultures, leave "trails" about various aspects of their lives...

"Social systems are systems of communication"

(Niklas Luhman)

By analyzing the "trails" (the results of social systems structural coupling) we can analyze the social system itself:politics, culture, media, business, technology, science ...

Outline



Introduction

Web 2.0, Semantic Web, Web 3.0

Network science

Data collection

Data analysis

Case study - CROSBI



The "new" science of networks





source: The New York Times (April 3rd, 1933., p. 17).



source: An Attraction Network in a Fourth Grade Class (Moreno, Who shall survive?, 1934).

Political and financial networks





Mark Lombardi (1980s & 1990s)



Terrorist networks



11 of 54

Board of directors membership networks





source: http://theyrule.net



On-line social networks







source: Bill Cheswick http://www.cheswick.com/ches/map/gallery/index.html

Air traffic networks





source: Northwest Airlines WorldTraveler Magazine

Railway networks





source: TRTA, March 2003 - Tokyo rail map

Semantic networks





source: http://wordnet.princeton.edu/man/wnlicens.7WN

17 of 54



Gene networks



source: http://www.zaik.uni-koeln.de/bioinformatik/regulatorynets.html.en



Food chain networks



19 of 54



Where can we find networks?

social networking sites



- social networking sites
- forums, buletinboards, discussion groups, mailing lists



- social networking sites
- forums, buletinboards, discussion groups, mailing lists
- blogs, microblogs, vlogs



- social networking sites
- forums, buletinboards, discussion groups, mailing lists
- blogs, microblogs, vlogs
- wikis, semantic wikis



- social networking sites
- forums, buletinboards, discussion groups, mailing lists
- blogs, microblogs, vlogs
- wikis, semantic wikis
- social tagging



- social networking sites
- forums, buletinboards, discussion groups, mailing lists
- blogs, microblogs, vlogs
- wikis, semantic wikis
- social tagging
- RSS feeds, RDF repositories, metadata collections



- social networking sites
- forums, buletinboards, discussion groups, mailing lists
- blogs, microblogs, vlogs
- wikis, semantic wikis
- social tagging
- RSS feeds, RDF repositories, metadata collections...



Where can we find networks?

- social networking sites
- forums, buletinboards, discussion groups, mailing lists
- blogs, microblogs, vlogs
- wikis, semantic wikis
- social tagging
- RSS feeds, RDF repositories, metadata collections...

A social network can be found on any application where users communicate.



Where can we find networks?

- social networking sites
- forums, buletinboards, discussion groups, mailing lists
- blogs, microblogs, vlogs
- wikis, semantic wikis
- social tagging
- RSS feeds, RDF repositories, metadata collections...

A social network can be found on any application where users communicate.

A semantic network can be found in any specification where concepts are meaningfully connected.

Outline



Introduction

Web 2.0, Semantic Web, Web 3.0

Network science

Data collection Semistructured data XPath Regular expressions Ajax APIs

Data analysis

21 of 54 study - CROSBI


Common obstacles in data collection

Semi-structured data (HTML, XHTML, XML, RDF, OWL ...)



Common obstacles in data collection

Semi-structured data (HTML, XHTML, XML, RDF, OWL ...)



The need for pattern matching

TOP Administration (2009-11-26 13:02:10) - Markus Schatten mpare to previous (23339)

TOP Administration (2011-01-30 05:47:00) - Markus Schatten impare to previous (25144)

TOP Administration (2011-01-30 05:48:26) - Markus Schatten mpare to previous (44477)

22 of 54

Common obstacles in data collection



Semi-structured data (HTML, XHTML, XML, RDF, OWL ...)



The need for pattern matching

TOP Administration (2009-11-26 13:02:10) - Markus Schatten ompare to previous (23339)

TOP Administration (2011-01-30 05:47:00) - Markus Schatten impare to previous (25144)

TOP Administration (2011-01-30 05:48:26) - Markus Schatten mpare to previous (44477)



Semistructured data and XPath



- query language for semistructured data
- part of the W3C standard
- uses path expressions to navigate through documents
- has a standard library with over 100 useful functions Example:

/html/body/div[@id='background']/div/div[@id='mainBar']/div/p[last()-1]

XPather (Firefox add-on)



Pogosta vprašanja Alije FiŠ javna akuteta. Ustava FiŠ je javna akuteta. Ustava Ali je fakuteta samostojna ali FiŠ je samostojni javil visokosti FiŠ je samostojni javil visokosti

Ali je študij brezplačen?

Za redne študente šolnine ni, plačajo samo vpisne stroške.

Kje potekajo predavanja?

Predavanja potekajo v Novem mestu (Kulturni center Janeza (Gimnazija Nova Gorica, Delpinova 9 in Dijaški dom) v Ljublja

Kje se nahaja fakulteta?

	😢 XPather Browser					×	
	XPath v	und']/div/div[@id='mainBar'		ar']/div[2]	r']/div[2]/p/strong		?
and the second second	RegExp			Subst			
Summer 12	Matching M	Nodes (count:	10 from 10)			
ašanja	no f	ull XPath html/body/div[@id='backgro	und'1/div/	div[@id='r	nainBar'l/di	E\$
i zasebna fakulteta?	2 /r 3 /r	ntml/body/div[ntml/body/div[@id='backgro @id='backgro	und']/div/ und']/div/	div[@id='r div[@id='r	nainBar']/di nainBar']/di	
ilteta. Ustanovljena je t	4 /r	tml/body/div[@id='backgro	und']/div/	div[@id='r	nainBar']/di	
lije v Novem mestu (OdU	5 /r	tml/body/div[@id='backgro	und']/div/	div[@id='r	nainBar']/di	
ımostojna ali del katere	Content of	the selected	@id=!backom nodes	und 1/div/	divi@id='r	nainRar'1/di	4
javni visokošolski zavod	Text Inne	er HTML/XML	Web Clipping	XPaths	Info		
lačen?	Ali je Fl	Ś javna ali zas	sebna fakultet	a7			٦
) šolnine ni, plačajo sam	<u> </u>	,					
Javanja?							
ajo v Novem mestu (Kult							- 1
Sorica Delninova 9 in Di							

Regular expressions



Useful for extracting various formatted data and searching for keywords

Example (date matching)

 $27-9-2001 \qquad [0-9] \{1,2\} [0-9] \{1,2\} [1-2] [0-9] \{3\}$



 Some sites use Ajax to generate data on page load or user interaction (e.g. mouse click, mouse over etc.)



- Some sites use Ajax to generate data on page load or user interaction (e.g. mouse click, mouse over etc.)
- In combination with obfuscated JavaScript code this is hard to scrape!
- Current solutions:
 - Browser scripting (can be slow and cumbersome) e.g. with selenium



- Some sites use Ajax to generate data on page load or user interaction (e.g. mouse click, mouse over etc.)
- In combination with obfuscated JavaScript code this is hard to scrape!
- Current solutions:
 - Browser scripting (can be slow and cumbersome) e.g. with selenium
 - □ Headless browser e.g. mechanize, spidermonkey, HtmlUnit



- Some sites use Ajax to generate data on page load or user interaction (e.g. mouse click, mouse over etc.)
- In combination with obfuscated JavaScript code this is hard to scrape!
- Current solutions:
 - Browser scripting (can be slow and cumbersome) e.g. with selenium
 - □ Headless browser e.g. mechanize, spidermonkey, HtmlUnit





 Most popular social web sites have an open API that allow you to connect your application to their site (Facebook, Twitter, Wikipedia, ...)



- Most popular social web sites have an open API that allow you to connect your application to their site (Facebook, Twitter, Wikipedia, ...)
- Still such APIs often have drawbacks for network data collection:



- Most popular social web sites have an open API that allow you to connect your application to their site (Facebook, Twitter, Wikipedia, ...)
- Still such APIs often have drawbacks for network data collection:
 - □ Limited access rights (e.g. Facebook)



- Most popular social web sites have an open API that allow you to connect your application to their site (Facebook, Twitter, Wikipedia, ...)
- Still such APIs often have drawbacks for network data collection:
 - □ Limited access rights (e.g. Facebook)
 - □ Limited queries per time interval (e.g. Twitter) etc.

Outline



Introduction

Web 2.0, Semantic Web, Web 3.0

Network science

Data collection

Data analysis Network matrices Folksonomy model Logic programming

Case study - CROSBI

28 of 54



Networks are most often analyzed using graph theory:

Definition

A graph \mathcal{G} is the pair $(\mathcal{N}, \mathcal{E})$ whereby \mathcal{N} represents the set of verticles or nodes, and $\mathcal{E} \subseteq \mathcal{N} \times \mathcal{N}$ the set of edges connecting pairs from \mathcal{N} .



Networks are most often analyzed using graph theory:

Definition

A graph \mathcal{G} is the pair $(\mathcal{N}, \mathcal{E})$ whereby \mathcal{N} represents the set of verticles or nodes, and $\mathcal{E} \subseteq \mathcal{N} \times \mathcal{N}$ the set of edges connecting pairs from \mathcal{N} .

Definition

A directed graph or digraph \mathcal{G} is the pair $(\mathcal{N}, \mathcal{A})$, whereby \mathcal{N} represents the set of nodes, and $\mathcal{A} \subseteq \mathcal{N} \times \mathcal{N}$ the set of ordered pairs of elements from \mathcal{N} that represent the set of graph arcs.



Networks are most often analyzed using graph theory:

Definition

A graph \mathcal{G} is the pair $(\mathcal{N}, \mathcal{E})$ whereby \mathcal{N} represents the set of verticles or nodes, and $\mathcal{E} \subseteq \mathcal{N} \times \mathcal{N}$ the set of edges connecting pairs from \mathcal{N} .

Definition

A directed graph or digraph \mathcal{G} is the pair $(\mathcal{N}, \mathcal{A})$, whereby \mathcal{N} represents the set of nodes, and $\mathcal{A} \subseteq \mathcal{N} \times \mathcal{N}$ the set of ordered pairs of elements from \mathcal{N} that represent the set of graph arcs.

Definition

A valued or weighted graph $\mathcal{G}_{\mathcal{V}}$ is the triple $(\mathcal{N}, \mathcal{A}, \mathcal{V})$ whereby \mathcal{N} represents the set of nodes or verticles, $\mathcal{A} \subseteq \mathcal{N} \times \mathcal{N}$ the set of pairs of elements from \mathcal{N} that represent the set of graph arcs, and $\mathcal{V} : \mathcal{N} \to \mathbb{R}$ a function that attaches values or weights to nodes.



Definition

A bipartite graph $\mathcal{G} = (\mathcal{N}, \mathcal{E})$ is a graph for which set of nodes there exists an partition $\mathcal{N} = \{X, Y\}$, such that every arc has one end in *X*, a and the other in *Y*.



Definition

A bipartite graph $\mathcal{G} = (\mathcal{N}, \mathcal{E})$ is a graph for which set of nodes there exists an partition $\mathcal{N} = \{X, Y\}$, such that every arc has one end in *X*, a and the other in *Y*.

Definition

A multigraph $\mathcal{G} = (\mathcal{N}, \mathcal{E})$ is a graph in which \mathcal{E} is a multiset, e.g. there can exist multiple edges between two nodes.

Matrix representation



Definition

Let G be a graph defined with the set of nodes $\{n_1, n_2, ..., n_m\}$ and edges $\{e_1, e_2, ..., e_l\}$. For every i, j ($1 \le i \le m$ and $1 \le j \le m$) we define

$$a_{ij} = \begin{cases} 1, & \text{if there is an edge between nodes } n_i \text{ and } n_j \\ 0, & \text{otherwise} \end{cases}$$

Matrix $A = [a_{ij}]$ is then the adjacency matrix of graph G. The matrix i symmetric since if there is an edge between nodes n_i and n_j then clearly there is also an edge between n_j and n_i . Thus $A = [a_{ij}] = [a_{ji}]$.

Simple graph





Digraph





33 of 54

Digraph





Weighted graph





	Varaždin	Zagreb	Rijeka	Split
Varaždin	0	98	280	463
Zagreb	98	0	175	365
Rijeka	280	175	0	405
Split	463	365	405	0

Bipartite graph





			۲			
	1	1	0	0	0	0
er	0	1	1	0	0	0
rade	0	0	0	1	1	0
	0	0	0	1	1	0
nder	0	0	0	1	1	1

36 of 54

Multigraph







Matrices can be transposed, multiplied ...



Matrices can be transposed, multiplied ... Examples:

 \mathbb{P} - process flow matrix ($p \times p, p$ - number of processes)



Matrices can be transposed, multiplied ... Examples:

 \mathbb{P} - process flow matrix ($p \times p, p$ - number of processes)

 \mathbb{W} - co-worker matrix ($w \times w$, w - number of workers)



Matrices can be transposed, multiplied ... Examples:

- $\mathbb P$ process flow matrix ($\rho \times \rho, \, \rho$ number of processes)
- \mathbb{W} co-worker matrix ($w \times w$, w number of workers)
- $\mathbb O$ process responsibility matrix ($p \times z$)



Matrices can be transposed, multiplied ... Examples:

- \mathbb{P} process flow matrix ($p \times p, p$ number of processes)
- \mathbb{W} co-worker matrix ($w \times w$, w number of workers)
- $\mathbb O$ process responsibility matrix (p imes z)
- WW co-workers of co-workers



Matrices can be transposed, multiplied ... Examples:

 \mathbb{P} - process flow matrix ($p \times p, p$ - number of processes)

 \mathbb{W} - co-worker matrix ($w \times w$, w - number of workers)

 $\mathbb O$ - process responsibility matrix ($p \times z$)

WW - co-workers of co-workers

 $\mathbb{P}^{\mathcal{T}}$ - inversed process flow



Matrices can be transposed, multiplied ... Examples:

 \mathbb{P} - process flow matrix ($p \times p, p$ - number of processes)

 \mathbb{W} - co-worker matrix ($w \times w$, w - number of workers)

 $\mathbb O$ - process responsibility matrix ($p \times z$)

WW - co-workers of co-workers

 $\mathbb{P}^{\mathcal{T}}$ - inversed process flow

 $\mathbb{O}^T\mathbb{O}$ - matrix of workers who are responsible for the same processes



Matrices can be transposed, multiplied ... Examples:

 \mathbb{P} - process flow matrix ($p \times p, p$ - number of processes)

 \mathbb{W} - co-worker matrix ($w \times w$, w - number of workers)

 $\mathbb O$ - process responsibility matrix ($p \times z$)

WW - co-workers of co-workers

 $\mathbb{P}^{\mathcal{T}}$ - inversed process flow

 $\mathbb{O}^T\mathbb{O}$ - matrix of workers who are responsible for the same processes

 $\mathbb{P}\mathbb{O}$ - matrix of workers who are responsible for the forthcoming processes



Matrices can be transposed, multiplied ... Examples:

 \mathbb{P} - process flow matrix ($p \times p, p$ - number of processes)

W - co-worker matrix ($w \times w$, w - number of workers)

 $\mathbb O$ - process responsibility matrix ($p \times z$)

WW - co-workers of co-workers

 $\mathbb{P}^{\mathcal{T}}$ - inversed process flow

 $\mathbb{O}^T\mathbb{O}$ - matrix of workers who are responsible for the same processes

 $\mathbb{P}\mathbb{O}$ - matrix of workers who are responsible for the forthcoming processes

 $\mathbb{P}^{T}\mathbb{O}$ - matrix of workers who are responsible for the previous processes
Matrix algebra and interpretation



Matrices can be transposed, multiplied ... Examples:

 \mathbb{P} - process flow matrix ($p \times p, p$ - number of processes)

 \mathbb{W} - co-worker matrix ($w \times w$, w - number of workers)

 $\mathbb O$ - process responsibility matrix ($p \times z$)

WW - co-workers of co-workers

 $\mathbb{P}^{\mathcal{T}}$ - inversed process flow

 $\mathbb{O}^T\mathbb{O}$ - matrix of workers who are responsible for the same processes

 $\mathbb{P}\mathbb{O}$ - matrix of workers who are responsible for the forthcoming processes

 $\mathbb{P}^{T}\mathbb{O}$ - matrix of workers who are responsible for the previous processes

• • •

38 of 54



A - actors



- A actors
- C concepts



- A actors
- C concepts
- I instances



- A actors
- C concepts
- I instances
- $T \subseteq A \times C \times I$ folksonomy



- A actors
- C concepts
- I instances
- $T \subseteq A \times C \times I$ folksonomy

A tripartite hypergraph which can be reduced into three bipartite graphs:



- A actors
- C concepts
- I instances
- $T \subseteq A \times C \times I$ folksonomy

A tripartite hypergraph which can be reduced into three bipartite graphs:

• AC - network of actors and concepts



- A actors
- C concepts
- I instances
- $T \subseteq A \times C \times I$ folksonomy

A tripartite hypergraph which can be reduced into three bipartite graphs:

- AC network of actors and concepts
- *AI* network of actors and instances



- A actors
- C concepts
- I instances

• $T \subseteq A \times C \times I$ - folksonomy

A tripartite hypergraph which can be reduced into three bipartite graphs:

- AC network of actors and concepts
- AI network of actors and instances
- CI network of concepts and instances



Delicious - analysis



40 of 54

Problems with matrix representation





Problems with matrix representation





Only the "1"-s contain information

41 of 54

Sparse matrix











If A is connected to B, then B is connected to A (undirected graphs)



- If A is connected to B, then B is connected to A (undirected graphs)
- How to find friends of friends?



- If A is connected to B, then B is connected to A (undirected graphs)
- How to find friends of friends?
- Is there a path from node A to node B? Which path is that?



- If A is connected to B, then B is connected to A (undirected graphs)
- How to find friends of friends?
- Is there a path from node A to node B? Which path is that?
- How to compute the previously outlined graphs?



- If A is connected to B, then B is connected to A (undirected graphs)
- How to find friends of friends?
- Is there a path from node A to node B? Which path is that?
- How to compute the previously outlined graphs?

Logic programming as a solution



Deductive languages like Prolog, *FLORA-2*

44 of 54

Logic programming as a solution



- Deductive languages like Prolog, *F*LORA-2
- Advantages:
 - Declarative problem solving ("on a higher level")
 - Very appropriate for expressing mathematical structures

Logic programming as a solution



- Deductive languages like Prolog, *FLORA-2*
- Advantages:
 - Declarative problem solving ("on a higher level")
 - Very appropriate for expressing mathematical structures
- Disadvantages:
 - Combinatoric explosion (use heuristics and predicate tabling)



If A is connected to B, then B is connected to A

link(A, B) :- link(B, A).

45 of 54



Friend of a friend

ff(A, B) :- link(A, C), link(C, B).

46 of 54



Is there a path from node A to node B?

```
path( A, B ) :- link( A, B ).
path( A, B ) :- link( A, C ), path( C, B ).
```

Which path is that?

whichpath(X, Y, [X, Y]) :- link(X, Y).
whichpath(X, Y, [X, Z | T]) : link(X, Z),
 whichpath(Z, Y, [Z | T]).

Outline



Introduction

Web 2.0, Semantic Web, Web 3.0

Network science

Data collection

Data analysis

Case study - CROSBI



Croatian Scientific Bibliography (http://bib.irb.hr)



- Croatian Scientific Bibliography (http://bib.irb.hr)
- Scientists enter data about their research activities by them selves which makes it a social application



- Croatian Scientific Bibliography (http://bib.irb.hr)
- Scientists enter data about their research activities by them selves which makes it a social application
- Research hypotheses:



- Croatian Scientific Bibliography (http://bib.irb.hr)
- Scientists enter data about their research activities by them selves which makes it a social application
- Research hypotheses:
 - analyze two social systems (the Croatian scientific community and the Croatian public) and see how the most important research concepts (as viewed by the scientific community) are interpreted by the public (does the public understand most important concepts from Croatian science in the last decade?).



- Croatian Scientific Bibliography (http://bib.irb.hr)
- Scientists enter data about their research activities by them selves which makes it a social application
- Research hypotheses:
 - 1. analyze two social systems (the Croatian scientific community and the Croatian public) and see how the most important research concepts (as viewed by the scientific community) are interpreted by the public (does the public understand most important concepts from Croatian science in the last decade?).
 - analyze how various scientific fields are conceptually and socially interrelated (do scientists do interdisciplinary research together if they have conceptually connected fields?)



Methodology

- Social systems are represented by their "trails": CROSBI as the Croatian scientific community, Croatian Wikipedia as the Croatian public
- One system can interpret concepts from another if the concept exists within its conceptual network (e.g. keyword or wiki page)

Data collection



- CROSBI data has been collected using Scrapy in November 2010 (a total of 285,234 entries has been collected)
- Wikipedia data has been collected using the Wikipedia API + Scrapy for parsing

Results



Bibliography



- Adamic, L.: Why networks are interesting to study, University of Michigan, School of Information, https://open.umich.edu/education/si/si508/fall2008
- Barratm M., Barthlemy, M., Vespignani, A.: Dynamical Processes on Complex Networks, Cambridge University Press, 2008.
- Carley, K.M.: Dynamic Network Analysis, Summary of the NRC workshop on social network modeling and analysis, Committee on Human Factors, National Research Council, 133145, Eds: Breiger, R., Carley, K.M., and Pattison, P., 2003.
- 4. Divjak, B., Lovrenčić, A.: Diskretna matematika s teorijom grafova, TIVA & FOI, 2005
- Krackhardt, David, and Carley, Kathleen M.: A PCANS Model of Structure in Organization, Proceedings of the 1998 International Symposium on Command and Control Research and Technology, Evidence Based Research, Vienna, VA, June 1998
- Mika, Peter: Ontologies are us: A unified model of social networks and semantics, Web Semantics: Science, Services and Agents on the World Wide Web 5(1), volume 5, 515, March 2007
- Newman, M., Barabsi, A.-L., Watts, D. j.: The Structure and Dynamics of Networks, Princeton University Press, 2006.
- 8. Various web sites (images, graphs)