Text Visualization

Maneesh Agrawala

CS 294-10: Visualization Spring 2011

Last Time: Graph Layout 2

Graphs and Trees

Graphs

Model relations among data *Nodes* and *edges*



Trees

Graphs with hierarchical structure Connected graph with N-1 edges Nodes as *parents* and *children*





v11 v12 v13 v14

Reverse some edges to remove cycles Assign nodes to hierarchy layers → Longest path layering Create dummy nodes to "fill in" missing layers Arrange nodes within layer, minimize edge crossings Route edges – layout splines if needed

v10

v8 v9

Layer 4

















Use radial tree layout for inner circle Mirror to outside Replace inner tree with hierarchical edge bundles



Summary

╶╠┈─╓╖──╔╗──╒╦╗─ Tree Layout

Indented / Node-Link / Enclosure / Layers How to address issues of scale?

Filtering and Focus + Context techniques

Graph Layout

Tree layout over spanning tree Hierarchical "Sugiyama" Layout Optimization (Force-Directed Layout) Attribute-Driven Layout



Final project

Design new visualization method

Pose problem, Implement creative solution

Deliverables

- Implementation of solution
- 8-12 page paper in format of conference paper submission
 2 design discussion presentations

Schedule

- Project proposal: 3/14
- Project presentation: 4/4
- Final paper and presentation: 5/3 1:30-3pm 6th floor Soda

Grading

- Groups of up to 3 people, graded individually
 Clearly report responsibilities of each member

Text Visualization

Why visualize text?

Why Visualize Text?

Understanding: get the "gist" of a document

Grouping: cluster for overview or classification

- Compare: compare document collections, or inspect evolution of collection over time
- Correlate: compare patterns in text to those in other data, e.g., correlate with social network

What is text data?

Documents

Articles, books and novels Computer programs E-mails, web pages, blogs Tags, comments

Collection of documents



Ú/čeur

9

Challenge: Visualize Dissertations

You have 20 years of university Ph.D. theses:

Text

- Year
- Dept.
- Author
- Advisor
- Committee

What questions might you want to answer? What visualizations might help?

A Concrete Example

What would help you gauge? Topics in document?

Relationship to other docs?

Supporting Asynchronous Collaboration for Interactive Visualization

Jeffrey Michael Heer

B.S. (University of California, Berkeley) 2001 M.S. (University of California, Berkeley) 2004

A dissertation submitted in partial satisfaction of the requirements for the degree of Doctor of Philosophy in Computer Science in the GRADUATE DIVISION

of the UNIVERSITY OF CALIFORNIA, BERKELEY









A Double Gulf of Evaluation

Many (most?) text visualizations do not represent text directly, they represent a model term statistics

- clusters
- 01000
- Can you interpret the visualization?
- How well does it convey the properties of the model?

Do you trust the model? How does the model enable us to reason about the text?

Lessons for Text Visualization

- Show (or provide access to) source text Let readers assess model Let readers use visualization as index into documents
- Find meaningful abstractions for grouping docs Are clusters interpretable?
- Where possible use text to represent text... but which terms are the most descriptive?

Topics

Text as Data Visualizing Document Content Evolving Documents Visualizing Conversation Document Collections

Text as Data

Words are (not) nominal?

High dimensional (10,000+) More than equality tests Words have meanings and relations

- Correlations: Hong Kong, San Francisco, Bay Area
- Order: April, February, January, June, March, May
- Membership: Tennis, Running, Swimming, Hiking, Piano
- Hierarchy, antonyms & synonyms, entities, ...

Text Processing Pipeline

Tokenization: segment text into terms

Special cases? e.g., "San Francisco", "L'ensemble", "U.S.A." Remove stop words? e.g., "a", "an", "the", "to", "be"?

Stemming: one means of normalizing terms

- Reduce terms to their "root"; Porter's algorithm for English e.g., automate(s), automatic, automation all map to automat
- For visualization, want to reverse stemming for labels

 Simple solution: map from stem to the most frequent word
- Result: ordered stream of terms

The Bag of Words Model

Ignore ordering relationships within the text

- A document ~ vector of term weights
- Each term corresponds to a dimension (10,000+) Each value represents the relevance

For example, simple term counts

Aggregate into a document x term matrix Document vector space model

Document x Term matrix

Each document is a vector of term weights Simplest weighting is to just count occurrences

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	157	73	0	0	0	0
Brutus	4	157	0	1	0	0
Caesar	232	227	0	2	1	1
Calpurnia	0	10	0	0	0	0
Cleopatra	57	0	0	0	0	0
mercy	2	0	3	5	5	1
worser	2	0	1	1	1	0

	WORDCOUNT
PREVIOUS WORD	NEXT WORD
the	(1988) (b) (b) (b) (b) (b) (b) (b) (b) (b) (b
	1999 Miller Miller Street St



Weaknesses of Tag Clouds

Sub-optimal visual encoding (size vs. position) Inaccurate size encoding (long words are bigger) May not facilitate comparison (unstable layout) Term frequency may not be meaningful Does not show the structure of the text

Keyword Weighting

 $\begin{array}{l} \textbf{Term Frequency} \\ tf_{td} = \mathsf{count}(t) \text{ in d} \\ \text{Can take log frequency: } \log(1 + tf_{td}) \\ \text{Can normalize to show proportion: } tf_{td} \ / \ \Sigma_t \ tf_{td} \end{array}$



Keyword Weighting

Term Frequency tf_{td} = count(t) in d

TF.IDF: Term Freq by Inverse Document Freq tf.idf_{td} = $log(1 + tf_{td}) \times log(N/df_t)$ df_t = # docs containing t; N = # of docs



Visualizing Document Content

osteoporosis prevention	Search Limit: \diamond 50 \diamond 100 \blacklozenge 250 \diamond 500 - 1			
research	Number of Clusters: $\bigcirc 3 \bigcirc 4 \spadesuit 5 \bigcirc 8 \bigcirc 1$			
Mode: TileBars				
Cluster Titles	Backup			
	FR88513-0157			
	AP: Groups Seek \$1 Billion a Year for Aging Research			
	SJMN: WOMEN'S HEALTH LEGISLATION PROPOSED CF			
	AP: Older Athletes Run For Science			
I water a second state	FR: Committee Meetings			
المحمدة ومحتك	FR: October Advisory Committees; Meetings			
	FR88120-0046			
	FR: Chronic Disease Burden and Prevention Models; Program			
	AP: Survey Says Experts Split on Diversion of Funds for AIDS			
	FR: Consolidated Delegations of Authority for Policy Developm			
	SJMN: RESEARCH FOR BREAST CANCER IS STUCK IN P			
RIR				























Glimpses of structure

Concordances show local, repeated structure But what about other types of patterns?

For example

Lexical: Syntactic: <A> at <Noun> <Verb> <Object>

Phrase Nets [van Ham et al]

- Look for specific linking patterns in the text: 'A and B', 'A at B', 'A of B', etc Could be output of regexp or parser
- Visualize extracted patterns in a node-link view Occurrences → Node size Pattern position → Edge direction





















Visualizing Revision History					
How to depict contributions over time?					
Example: Wikipedia history log					
Chocolate					
Legend: (cur) = difference with current version. (last) = difference with preceding version. M = minor edit					
(aa) (aai) . 1201.120.4 Mag 2003 . Duranching (meaten to do, rearrange ace alzo) (aa) (aai) . 11.59, 20.4 Mag 2003 . 2 Bundi (aa) (aai) . 11.52, 20.4 Mag 2003 . 2 Bundi (aai) (aai) . 11.52, 20.4 Mag 2003	on				





*	syn diff: ss	hconse	leis				
Diff	style: Side-by-side Y Enable syntax coloring						
				~			
Files	Changed:						
	. <u>sshconsole.is</u> : 1 change [<u>1</u>]						
/h	/nome/toddw/src/ssnconsole-read-only/content asnconsole.js						
	50 lines hidden (Equand)						
51	_term = new VT100(00, 24, "term");	51	_term = new VT100(00, 24, "term");				
62	//_term.debug_ = 1;	52	//_term.debug_ = 1;				
63	_term.curs_set(true, true, _term_box_element);	63	_term.curs_set(true, true, _term_box_element);				
54	_term.noecho();	54	_term.noecho[];				
55		55					
22	// Replace the go_getch_ function with our own, this is called	-	// Replace the go_getch_ function with our own, this is called				
20	// for every keypress that is passed through the terminal to the	60	// for every keypress that is passed through the terminal to the				
	// restricted VT100 character, sequence(s)		// resulted VTI00 character sequence(s)				
00	VT100.co getch = function() (90	VT100.co oetch = function() (
61	var vt = VT109, the vt ;	61	var vt = VT100, the vt ;				
62	if (vt undefined) (62	if (vt somevalue) (
63	return;	63	return;				
64)	64)				
65	<pre>var ch = vt.key_bufshift();</pre>	65	<pre>var ch = vt.key_bufshift();</pre>				
66	<pre>//dump("go_getch_:: ch: '" + ch + "'\n");</pre>						
67	if (ch undefined) {	66	if (ch undefined) {				
63	return;	67	return;				
69	· · · · · · · · · · · · · · · · · · ·	65					
74	11 (w.edu) as $0(equ) = 1)$ ($11 (vr.edu_{-} \text{ ss} (u.temp) = 1) ($				
		75	vt.mater(cr);				
72)	72)				
73	if (sah channel) (73	if (sah channel) (
74	_ssh_channel.send5tdin(ch);	74	_ssh_channel.send5tdin(ch);				
75)	75)				
76)	76)				
77		77					
78	<pre>var serverTextbox = document.getElementById("sshconsole_server_textbox");</pre>	78	<pre>var serverTextbox = document.getElementById('sshconsole_server_textb</pre>	ox");			
12	var connectionText;	0	var connectionText;				
	<pre>in (connectionText is window arguments[0] (connectionText; </pre>	-00	<pre>if (connectionText is visited armometal@l connectionText;</pre>				
82) else (82) else (
0	, (174 in	ts hidden [Example				
				×			







Visualizing Conversation

Visualizing Conversation

Many dimensions to consider:

Who (senders, receivers) What (the content of communication) When (temporal patterns)

Interesting cross-products: What x When → Topic "Zeitgeist" Who x Who → Social network Who x Who x What x When → Information flow











