CSE512 :: 4 Mar 2014 Text Visualization



Jason Chuang University of Washington

Text visualization: What-Why-How

What is text data?

Documents

Articles, books and novels E-mails, web pages, blogs

Text Snippets

Tweets, SMS messages Tags, comments, profiles

And More...

Computer programs, logs This slide! Collections of documents



Why visualize text?

Understanding – read a document
Summaries – get the "gist" of a document
Clustering – group together similar contents
Quantify – convert to numerical measures
Correlate – compare patterns in text to those in other data, e.g., correlate with social network

Example: Health Care Reform

Recent history Initiatives by President Clinton Overhaul by President Obama Text data News articles Speech transcriptions Legal documents

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

A Concrete Example

September 10, 2009

Obama's Health Care Speech to Congress

Following is the prepared text of President Obama's speech to Congress on the need to overhaul health care in the United States, as released by the White House.

Madame Speaker, Vice President Biden, Members of Congress, and the American people:

When I spoke here last winter, this nation was facing the worst economic crisis since the Great Depression. We were losing an average of 700,000 jobs per month. Credit was frozen. And our financial system was on the verge of collapse.

As any American who is still looking for work or a way to pay their bills will tell you, we are by no means out of the woods. A full and vibrant recovery is many months away. And I will not let up until those Americans who seek jobs can find them; until those businesses that seek capital and credit can thrive; until all responsible homeowners can stay in their homes. That is our ultimate goal. But thanks to the bold and decisive action we have taken since January, I can stand here with confidence and say that we have pulled this economy back from the brink.

I want to thank the members of this body for your efforts and your support in these last several months, and especially those who have taken the difficult votes that have put us on a path to recovery. I also want to thank the American people for their patience and resolve during this trying time for our nation.

But we did not come here just to clean up crises. We came to build a future. So tonight, I return to speak to all of yo

Tag Clouds: Word Count

President Obama's Health Care Speech to Congress [New York Times]



economix.blogs.nytimes.com/2009/09/09/obama-in-09-vs-clinton-in-93



economix.blogs.nytimes.com/2009/09/09/obama-in-09-vs-clinton-in-93

WordTree: Word Sequences

Visualizations : Word Tree President Obama's Address to Congress on Health Care

Search	Back Forward Start End Occurrence Order + Clicks Will Zoom +
52 hits	<complex-block><pre>interval in the interval interval</pre></complex-block>
	that we can act when it's hard.

Search i will	Back Forward Occurrence Order Clicks Will Zoom
12 hits	Intervention of the status quo as a solution . make sure that no government bureaucrat or insurance company bureaucrat gets between you and protect medicare . between you and in the weeks ahead .

neplace acrimony with civility, and gridlock with progress.
 de great things, and that here and now we will meet history's test.

we can -

- - i still believe that we can act when it's hard . - that we can act when it's hard .

still believe -

A Double Gulf of Evaluation

Many (most?) text visualizations do not represent the text directly. They represent the output of a language model (word counts, word sequences, etc.).

- 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?

Topics

Summarizing with Words Visualizing Themes in a Document Collection Quantifying Textual Content Performing Text Analysis

Summarize with Words

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

1. Tokenization

Segment text into terms. Remove stop words? *a, an, the, of, to, be* Numbers and symbols? *#gocard, @stanfordfball, Beat Cal!!!!!!!* Entities? San Francisco, O'Connor, U.S.A.

2. Stemming

Group together different forms of a word. Porter stemmer? visualization(s), visualize(s), visually \rightarrow visual Lemmatization? goes, went, gone \rightarrow go

3. Ordered list of terms

Tips: Tokenization and Stemming

Well-formed text to support stemming? *txt u l8r!*

Word meaning or entities? #berkeley → #berkelei

Reverse stems for presentation. Ha appl made programm cool? Has Apple made programmers cool?

Bag of Words Model

Ignore ordering relationships within the text

A document ≈ vector of term weights

- Each dimension corresponds to a term (10,000+)
- Each value represents the relevance
 For example, simple term counts

Aggregate into a document-term matrix

Document vector space model

Document-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 (Harris 2004)

			WORDCOUNT
PREVIOUS WORD			NEXT WORD
the	ofandtoain	thatitiswasi.foronyoutebraithasbyathermore this protection and the	D. M. Gan Star Strategy and s
CURRENT WORD			
FIND WORD:	BY RANK:	REQUESTED WORD: THE	86800 WORDS IN ARCHIVE

http://wordcount.org

Visualizations : Wordle of Sarah Palin RNC 9/3/2008 Speech

Creator: Anonymous Tags:

Edit Language Font Layout Color



Tag Clouds

Strengths

Can help with gisting and initial query formation.

Weaknesses

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

Term Frequency

tf_{td} = count(t) in d Can take log frequency: log(1 + tf_{td}) Can normalize to show proportion: tf_{td} / Σ_t tf_{td}



Partisan Words, 106th Congress, Abortion

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}) × log(N/df_t) df_t = # docs containing t; N = # of docs



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}) × log(N/df_t) df_t = # docs containing t; N = # of docs

G²: Probability of different word frequency $E_{1} = |d| \times (tf_{td} + tf_{t(C-d)}) / |C|$ $E_{2} = |C-d| \times (tf_{td} + tf_{t(C-d)}) / |C|$ $G^{2} = 2 \times (tf_{td} \log(tf_{td}/E_{1}) + tf_{t(C-d)} \log(tf_{t(C-d)}/E_{2}))$



Limitations of Frequency Statistics?

Typically focus on unigrams (single terms)

Often favors frequent (TF) or rare (IDF) terms Not clear that these provide best description

A "bag of words" ignores additional information Grammar / part-of-speech Position within document Recognizable entities

How do people describe text?

We asked 69 subjects (graduate students) to read and describe dissertation abstracts.

Students were given 3 documents in sequence; they then described the collection as a whole.

Students were matched to both familiar and unfamiliar topics; topical diversity within a collection was varied systematically.

[Chuang, Manning & Heer, 2012]





Term Commonness

$\log(tf_w) / \log(tf_{the})$

The normalized term frequency relative to the most frequent n-gram, e.g., the word "the".

Measured across an entire corpus or across the entire English language (using Google n-grams)





Scoring Terms with Freq, Grammar & Position



A fighter jet rain check

😴 🕂 🖉 http://www.boeing.com/Features/2010/09/bds_feat_morewater 🖒 🔍

A fighter jet rain check

Story and video by Chamila Jayaweera

Have you ever thought about what it takes to make sure that sea-based fighter jets stay dry?

When it comes to the F/A-18 Super Hornet, Boeing engineers in St. Louis use a special process called the Water Check Test to rule out areas where moisture could seep into the aircraft and its electronics suite.

Program experts douse the jet with simulated rain at a 15-inch-per-hour rate for about 20 minutes inside an enormous hangar in St. Louis.

"Our ultimate customers are U.S. Navy fighter pilots, and we want to ensure their safety in flight and on the ground, and water-tight integrity of the aircraft also



CHAMILA JAYAWEERA/BOEING

The Water Check team rolls in a large metal frame, which they affectionately call their "spray tree," over a Super Hornet inside a St. Louis hangar.

helps increase their effectiveness," said Boeing's Rich Baxter, F/A-18 Super Hornet final assembly manager.

To find out moreabout how the process works and watch the action unfold, click above to see the video story.


G²

Regression Model

fighter F/A Hornet Super Boeing -18 rain St. jet Louis 15-inch-per-hour douse hangar water-tight Check Baxter sea-based aircraft Rich seep click Navy sure Water moisture watch enormous stay

want.

Super Hornet F/A -18 fighter jet **Boeing engineers** special process rain check electronics suite Program experts simulated rain ultimate customers enormous hangar water-tight integrity Rich Baxter 15-inch-per-hour rate video story aircraft U.S. Navy fighter pilots Super Hornet final assembly manager U.S. Navyfighter tighter pilot sea-based tighter

Yelp: Review Spotlight [Yatani 2011]



Yelp: Review Spotlight [Yatani 2011]



Tips: Descriptive Keyphrases

Understand the limitations of your language model.

- Bag of words:
 - Easy to compute Single words Loss of word ordering

Select appropriate model and visualization Generate longer, more meaningful phrases Adjective-noun word pairs for reviews Show keyphrases within source text

Visualize Themes in a Document Collection

Topical Analysis

Large document collections Too large to manually read the source documents Deeper analysis than the most common theme Statistical topic modeling Analysis of word relationships Extract *latent topics* belonging to the documents

Statistical Topic Modeling



Statistical Topic Modeling



Computational Linguistics

[Hall et al. 2008]

ACL conferences and journals

14,000 papers over 40 years



Computational Linguistics

[Hall et al. 2008]



NIH Grants & Funding Agencies

Biomedical research

110,000 NIH grant awards 220,000 MEDLINE journal articles

Statistical topic modeling

700 latent topics

Clustering & correlations

Project areas & Funding institutes Trends in research funding Changes in research topics



NIH Grants & Funding Agencies

Biomedical research

110,000 NIH grant awards 220,000 MEDLINE journal article

Generate 700 latent topics Remove 15% "nonsensical" topics Modify vocabulary (phrases, acronyms) Extensive parameter search Expert validation and topic curation



AM

HD

RI

List of Words

Anaphora Resolution Automata Biomedical Call Routing Categorial Grammar Centering* Classical MT Classification/Tagging Comp. Phonology Comp. Semantics*

resolution anaphora pronoun discourse antecedent pronouns coreference string state set finite context rule algorithm strings language symbol medical protein gene biomedical wkh abstracts medline patient clinica call caller routing calls destination vietnamese routed router destinatio proof formula graph logic calculus axioms axiom theorem proofs lamb centering cb discourse cf utterance center utterances theory coherence japanese method case sentence analysis english dictionary figure japan features data corpus set feature table word tag al test vowel phonological syllable phoneme stress phonetic phonology prom semantic logical semantics john sentence interpretation scope logic for

[Hall et al. 2008]

<u>Topic Words</u>: vegf angiogenesis vascular_endothelial_growth_factor angiogenic er antiangiogenic anti_angiogenic vegf_a tumor_angiogenesis vegfr2 growth signaling <u>Title Words</u>: angiogenesis, vegf, vascular_endothelial_growth_factor, angiogenic, t neovascularization, angiopoietin, signaling, vegfr, vascular, human <u>Phrases</u>: vascular_endothelial_growth_factor vegf, vegf angiogenesis, vegf recepto

Termite | Topic Model Visualization





SpicyNodes: Radial Layout Authoring for the General Public









online



online



online



online



online



online



online



online



online



online



Termite | Topic Model Visualization





Filtering: What words to show?

Frequent words are not necessarily discriminative

data, visualization, information, visual, techniques, users, visualizations, ...

Saliency Score for word *w* based on frequency and distinctiveness

saliency(w) = frequency(w) × distinctiveness(w)

Distinctiveness

Knowing a word w, how much does it tell us about a topic?

distinctiveness(w) = KL(P(T|w) || P(T))

P(T|w) = generating topic T for a given word w P(T) = generating topic T for a randomly-selected word in the corpus

Seriation: How to show the words?





Seriation: How to show the words?

Clustering of related words large, node, networks, social, link, diagrams, online, dataset, communities, ...

Preservation of reading order

online communities, social networks, node link diagrams, large datasets, ...

Word similarity matrix (asymmetric)

document co-occurrence sentence co-occurrence collocation (word transition probability)

Text seriation

based on bond energy algorithm accepts asymmetric similarity matrices aware of salient terms early termination



Seriation: How to show the words?

Clustering of related words large, node, networks, social, link, diagrams, online, dataset, communities, ...

Preservation of reading order

online communities, social networks, node link diagrams, large datasets, ...

Word similarity matrix (asymmetric)

document co-occurrence sentence co-occurrence collocation (word transition probability)

Text seriation

based on bond energy algorithm accepts asymmetric similarity matrices aware of salient terms early termination







Quantify Textual Content

London Riot 2012

Select a rumour to see how unsubstantiated claims are spread on Twitter before being confirmed or denied



Rioters attack London zoo and release animals



Rioters cook their own food in McDonald's



Police 'beat a 16-year-old girl'



London Eye set on fire



Rioters attack a children's hospital in Birmingham



Army deployed in Bank



Miss Selfridge set on fire
How riot rumors spread?

Our last challenge was to classify each tweet according to a 'common senseunderstanding' of its main role as a communicative act. Did it support, oppose, query or comment on a rumour? In addition to an algorithmic analysis by our academic partners, each tweet was independently coded by three sociology PhD students in order to enable us to check for reliability. All the results were then subject to final review for quality assurance purposes. These categories could then be used to colour code each tweet so that readers get an overall picture of what direction the dialogue is taking.



TileBars [Hearst]

LIWC: Linguistic Inquiry & Word Count

Psychological Processes

Social processes	social	Mate, talk, they, child	455	
Family	family	Daughter, husband, aunt	64	.87
Friends	friend	Buddy, friend, neighbor	37	.70
Humans	human	Adult, baby, boy	61	
Affective processes	affect	Happy, cried, abandon	915	
Positive emotion	posemo	Love, nice, sweet	406	.41
Negative emotion	negemo	Hurt, ugly, nasty	499	.31
Anxiety	anx	Worried, fearful, nervous	91	.38
Anger	anger	Hate, kill, annoyed	184	.22
Sadness	sad	Crying, grief, sad	101	.07
Cognitive processes	cogmech	cause, know, ought	730	
Insight	insight	think, know, consider	195	
Causation	cause	because, effect, hence	108	.44
			-	

Visual Thesaurus [ThinkMap]



comments



Named Entity Recognition

Identify and classify named entities in text: John Smith → PERSON Soviet Union → COUNTRY 353 Serra St → ADDRESS (555) 721-4312 → PHONE NUMBER

Entity relations: how do the entities relate? Simple approach: do they co-occur in a small window of text?

🎕 List View



х

Doc. Similarity & Clustering

In vector model, compute distance among docs

- For TF.IDF, typically cosine distance
- Similarity measure can be used to cluster

Topic modeling approaches

- Assume documents are a mixture of topics
- Topics are (roughly) a set of co-occurring terms
- Latent Semantic Analysis (LSA): reduce term matrix
- Latent Dirichlet Allocation (LDA): statistical model

Parallel Tag Clouds [Collins et al 09]

affidavit adjourned abuse appeal ballo' accused about abuse adverted allocatur adequate agency bankruptcy bargaining argument banc agency alia alia argued analysis aliens affirmed benefit affirmed assistant brief asked agency's anent black annuity barge allocution antitrust appropriate attached aid capital coal appellee boat antidumping appellant arbitration app called cargo bindina authority asylum cars candidate arbitration application brief charter argued asbestos appeal appellant's bargaining broker asbestos court case art cited circuit closure brief assets coverage appellee argument conspiracy board commenced bankruptcy cited collateral broadcast defendant COLAS class damages because complaint crack cable asseveration claim believe contended before death defendant copy commerce capricious defendant's denied composition below court copyright benefit debtor fisciplinar carrier compounds boat coal disability drilling denied enough crime county distribution dba competition date bottlers construction determine brief disability declared court conspiracy fire district costs denial defendant disfavor contract death district contention gang class data deportation disenfranchised aas drug commonwealth disposition get doc discrimination autortice habeas gun discretion defendant court's employees decision evidence emissions doctrine context foreign disenfranchised del had crack description creditors homestead farm disposition employees filed firearm estoppel fraud exemption debtor dozer device ensued harassing indemnity district decisional firearm grams examination explanatory error disclosed election have iniury exercise follows errs forthwith had embodied facilities event denied electors instant grievance help factfinding injunction gas fiduciary except equivalent disclosed her furnished ferritin insurance immunities hazard her inter hereby further fear have quidelines dispensed jack gas interpretation his inequitable his him fish insurance interna grazing he here intervenor here distribution jurists impair ivorv infringement judgment incarcerative law habitat his keeplock district labor inasmuch job inmates jail hardship inference liability license insurance iudicata invalid drug judgment job marks jury his inter error law lien interest material jury fact magistrate judge immigration invention memoranda loan migrant medical millions iurisdiction magistrate's limned urisdiction marihuana nevertheless from inventor mitigation petitioner kilogram methamphetamine his land lst legislation maritime nonstatutory layer interlocutory opinion pipelines omitter lawyer months may might ioined mitigation motion liability ordinance preceding oral more might office legal plaintiff native means negligence office mortgage payday promulgated majority more opinion novo order lung plaintiff's merchandise nre panel proposed market phase magistrate plausible one pain method principal paupers persuasive ostrich panel preceding oil postconviction point notes out material rate noninfringement qualified plaintiff quotation parish persecution para quantity plaintiff's rescript merits our said regulations racketeering miner's plaintiff's precedential police racial patent petition parent platform reversed say reinsurance pension prisoner mining rehearing record search patentee policy pneumoconiosis section plaintiff respect prisoner sav see reprinted opinion recovery police sentence provided res sentence security she product some plan pulmonary oral ref'd sexual public submittee rulemaking suggested protein sheriff see pursuant plenary suit suspended order pursuant supra removed she section reissue shareholders retard review students supra recommendation policy retiremen see think shares pneumoconiosis rigging subd trial specie tit tentative service sterile search said seaman turtle tit present provision suitable stock than unanimous servitude tusks testified shipper sentence signal recognized subway vesse thought town process stated skill sitting testimony unfavorable tribal summation specific trialworthy section published told told vote tariff suit summary tribe trial pulmonary unpublished structure vessel trade unanimous settlement usury voters transmission tribal want unanimous surface vacated union syrup vesse value viz sentence union what white under verdict water vehicle upon vaccination waybill whom which writ wrote vol warrant work zone veterans waste where would without would First Fifth Seventh DC Second Third Fourth Sixth Eighth Ninth Tenth Eleventh Federal

ThemeRiver [Havre et al 99]



TIARA [Wei et al. 09]





THE WORDS THAT WERE USED

The 2007 State of the Union Address

Over the years, President Bush's State of the Union address has averaged almost 5,000 words each, meaning the the President has delivered over 34,000 words. Some words appear frequently while others appear only sporadically. Use the tools below to analyze what Mr. Bush has said.



The word in context

Next Instance of 'Tax'

I believe in local control of schools. We should not, and we will not, run public schools from Washington, D.C. Yet when the federal government spends **TAX** dollars, we must insist on results. Children should be tested on basic reading and math skills every year between grades three and eight. Measuring is the only way to know whether all our children are learning. And I want to know, because I refuse to leave any child behind in America.

-- 2001 (Paragraph 14 of 73)

New York Times



* As a newly elected president, Mr. Bush did not deliver a formal State of the Union address in 2001. His Feb. 27 speech to a joint session of Congress was analogous to the State of the Union, but without the title.

Concordance

What is the common local context of a term?

A Concordance	- Larkin.	Concordance					X
<u>File T</u> ext <u>S</u> earch	<u>E</u> dit <u>H</u> ea	adwords Conte <u>x</u> ts <u>V</u> iew T <u>o</u> o	ols Hel <u>p</u>				
🛇 # 🖻 🖥 é	🗿 👗 l		P	0			
Headword	No. 🔨		Context	Word	Context	Reference	
HEAR	15	That	my own	heart	drifts and cries, having no	Deep Analysis	Ce
HEARD	9	By the sho	out of the	heart	continually at work	And the wave	R.
HEARING	7	Nothing to adapt the s	kill of the	heart	to, skill	And the wave	8
HEARS	3	The tread, the beat of it, it is	my own	heart		Träumerei	
HEARSE	1	Because I follow it to	my own	heart		Many famous	Г
HEART	25		My	heart	is ticking like the sun:	lam washed u	5
HEART'S	2	Tł	ne vague	heart	sharpened to a candid co	The March Pa:	in the second s
HEART-SHAPED	1 👝	Cor	ntract my	heart	by looking out of date.	Lines on a Yo 📄	ligr
HEARTH	1	H	laving no	heart	to put aside the theft	Home is so Sa	De la
HEARTS	7	And the boy p	uking his	heart	out in the Gents	Essential Beau	
HEARTY	1	A harbou	ur for the	heart	against distress.	Bridge for the	
HEAT	6	These I would ch	ioose my	heart	to lead	After-Dinner F	
HEAT-HAZE	1	Time in his little ciner	na of the	heart		Time and Space	[] a
HEATH	1	This	petrified	heart	has taken,	A Stone Churc	×
HEATS	1	How should they sweep the g	irl clean	heart		l see a girl dra	
HEAVE	1	Hands	s that the	heart	can govern	Heaviest of flo	P
HEAVEN	4		For the	heart	to be loveless, and as col	Dawn	
HEAVEN-HOLDING	1	With the ungu	essed-at	heart	riding	One man walk	17
HEAVIER-THAN	1	If hands could t	free you,	heart		If hands could	9
HEAVIEST	2 🗸	That overf	lows the	heart		Pour away the 💙	a
<	3	<					
Worde	Tokens	At word Deleted lines	Word sort		Context sort		
7010	27070		Ass alshe	(atriage)	Aco coourses a	ardar	
7318	37070	2990 1 [24]	Asc alpha	(string)	Asc occurrence of	praer	



Down the Rabbit-Hole

lan pictures

hurwalk

ifenktier

afteen Svind

Down the Rabbit-Hole

Hide text Show concordance

Show thesaurus lookup Show story ine

ine un for the sole of the sol

regular Othe consider regular Othe consider ashamed construction bed ashamed construction bed

May 20ice) levau'v&urptisear&are, hudles. made indee interrupted : a lightly is .course exclaimed and most a Burganationidam and add donë begare, ethought user of anxigHalled inns þýlfser whisteered alice Rectastone: spinker

externation and a subsection of the subsection o

> extraordinary Cushionveritesentra Retratepted mad singing

Show only KWIC index (Key Word In Context)

Down the Rabbit Hole Rabbit with pink eyes ran close by h which it so VERY much out of the way to hear the Rabbit say to but when the Rabbit actually TOOK A WATCH OUT OF ITS WAIST COATbefore seen a rabbit with either a waistcoat-pocket, or a watch to down a large rabbit-hole under the hedge.

The subbt-hole went straight on like a tunnel for some way, was another long passage, and the White Rabbit was still in corner, but the Rabbit was no longer to be seen: she found it was the White Rabbit returning, splendidy dressed, with a that she was ready to ask help of any one, so, when the Rabbit sign." The Rabbit started violently, dropped the white kid supprised to see that, she had put on one of the Rabbit's ittle be self. The Rabbit Sends in a Little Bill.

It was the White Rabbit, trotting slowly back again, and Yery soon the Rabbit noticed Alice, as she went hunting about, messages for a rabbit 1 suppose Dinah'll be sending me on by mice and rabbits. I almost wish I hadn't gone down that abbit-hole--and yet--and yet--it's rather curious, you know, stairs. Alice knew it was the Rabbit coming to lock for her, and was now about a thousand times as large as the Rabbit, and had no Presently the Rabbit came up to the door, and tried to open it; lancied she heard the Rabbit just under the window, she suddenly 'Digging for apples, indeed!' said the Rabbit angrily. 'Here! began moving about again, and Alice heard the Rabbill say, "A recognised the White Rabbit: it was taking in a hurried nervous She was waking by the White Rabbit, who was peeping anxiously. 'Hush' Hush' said the Rabbit in a low, hurried tone. He Did you say "What a pity!"?" the Rabbit asked. "She boxed the Queen's ears--' the Rabbt began. Alice gave a ittle scream of laughter. "Oh, hush!" the Rabbit whispered in a she first saw the White Rabbit. She was a little nervous about each side to guard him; and near the King was the White Rabbit, she stopped hastiy, for the White Rabbit cried out, "Silence in On this the White Rabbit blew three blasts on the trumpet, and "Not yet, not yet!" the Rabbit hastily interrupted. "There's

"Call the first witness," said the King; and the White Rabbit

6 X

if love be rough with you, be rough with love.
if love be blind, love cannot hit the mark.
if love be blind, it best agrees with night.



WordTree (Wattenberg et al)



Filter infrequent runs

ling

art .		h
3		
8		
	The second	
Land /		
lora ┥		
	Contraction of the Contraction o	
	and the second sec	
-		
-		
Test.		
810		
-		
14		
	and the second sec	
the state of the s		
childree	d 10	
-		
bank a		
IDS	1	
100		
- Take	A second s	
harme.		
ECHINF .		
de-		
	have been been been been been been been be	
	and the second sec	
1		
and a second		
4 · · · ·		
2		
-		
£		
F		
Ì.		



Recurrent themes in speech



search

 \mathbf{v}

many eyes

explore visualizations data sets comments topic hubs

participate create visualization upload data set create topic hub register

learn more

quick start visualization types data format & style about Many Eyes FAQ blog

contact Us contact

report a bug

legal terms of use

Popular Dataset Tags 2007 2008 bible blog books CENSUS crime education eharmony election energy food health inauguration internet ireland literature lyrics media music network obama

people politics population

president prices religion

currently showing





Endings and Job Progressions

Word Transition Probability

Glimpses of structure

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

For example

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

Phrase Nets [van Ham et al 2009]

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



Node Grouping







eye lay glancing









Text Visualization Summary

High Dimensionality

Where possible use **text to represent text**... ... which terms are the most descriptive?

Context & Semantics

Provide **relevant context** to aid understanding. Show (or provide access to) the **source text**.

Modeling Abstraction

Determine your **analysis task**.

Understand abstraction of your language models.

Match analysis task with appropriate tools and models.

Perform Text Analysis

Information Retrieval

Search for documents Match query string with documents

Google scholar acronym resolution Search Advanced Scholar Search	
Scholar Articles and patents anytime include citations Create email alert	
A supervised learning approach to acronym identification D Nadeau, P Turney - The Eighteenth Canadian, 2005 - nparc.cisti-icist.nrc-cnrc.gc.ca Recently the fields of Genetics and Medicine have become especially interested in acronym resolution (Pustejovsky et al., 2001, Yu et al. 2002) Pustejovsky et al.'s acronym resolution technique searches for definitions of acronyms within noun phrases <u>Cited by 48</u> - <u>Related articles</u> - <u>All 16 versions</u>	[PDF] from nrc-cnrc.
Biomedical term mapping databases JD Wren, JT Chang, J Pustejovsky Nucleic acids, 2005 - Oxford Univ Press the prevalence of polynyms, or acronyms with multiple definitions. An important part of any high-throughput effort to tie experimental findings to published knowledge within the scientific literature involves acronym resolution <u>Cited by 41</u> - <u>Related articles</u> - <u>All 22 versions</u>	[HTML] from nih.gov Find it@Stanford
Anthropogenic climate change over the Mediterranean region simulated by a global variable resolution	Find it@Stanford
[TIODE] AL Gibelin Climate Dynamics, 2003 - Springer The long simulations CC and CS are split into two 30-year datasets CC1 and CS1 for the period 1960–1989 and CC2 and CS2 for the period 2070–2099 Full name Acronym Resolution Period Coupled Coupled control CC T63 1950–2099 Yes <u>Cited by 197 - Related articles - BL Direct</u> - <u>All 5 versions</u>	
Metaphrase: an aid to the clinical conceptualization and formalization of patient problems in healthcare enterprises. MS Tuttle, NE Olson, KD Keck, WG Cole Methods of information, 1998 - ukpmc.ac.uk Title not supplied (PMID: 10566483). Concept definition and manipulation are supported through	

	Tool	ls T <u>a</u> ble <u>W</u> indow <u>H</u> elp Ado	be PDF Acrobat <u>C</u> omments	
	ABC	Spelling and Grammar F7	100% 👻 🙆 🚉 Keau 📮 : 🚣 Normal 📼	
	í,	<u>R</u> esearch Alt+Click	bw + 😥 💫 🛷 + 💫 + 🦕 🥸 + 🍋 📄	
		Language •		
		<u>W</u> ord Count		
l		Speech	1	
		Shared Wor <u>k</u> space		
		<u>Track Changes</u> Ctrl+Shift+E		
		L <u>e</u> tters and Mailings	Track Changes Icons appear once you	
		<u>C</u> ustomize	select "Track Changes" from the "Tools	
		Options	Menu	

This is a test document to demonstrate the use of tracking changes. The characters in black font represent the original document while the characters in red font represent the changes which are being tracked.



Visualizing Revision History

How to depict contributions over time? Example: Wikipedia history log

Chocolate

Revision history

Legend: (cur) = difference with current version, (last) = difference with preceding version, M = minor edit

- (cur) (last) . . <u>12:01, 20 Aug 2003</u> . . <u>Dysprosia</u> (neaten to do, rearrange see also)
- (cur) (last) . . <u>11:59, 20 Aug 2003</u> . . <u>Patrick</u>
- (cur) (last) . . . <u>11:52, 20 Aug 2003</u> . . . <u>81.203.98.109</u>
- (cur) (last) . . M <u>18:36, 6 Aug 2003</u> . . <u>Manika</u> (corrected spelling)
- (cur) (last) . . <u>18:32, 6 Aug 2003</u> . . <u>Daniel Quinlan</u> (removing obscure heraldry information, belongs on [[heraldry]] if anywhere)
- (cur) (last) . . 15:21, 6 Aug 2003 . . Rmhermen
- (cur) (last) . . <u>15:08, 6 Aug 2003</u> . . <u>Cyp</u> (Chocolate often has odd shapes.)
- (cur) (last) . . <u>19:14, 3 Aug 2003</u> . . <u>Daniel C. Boyer</u> ("chocolate" as shade of gules in heraldry)
- (cur) (last) . . M <u>02:00, 30 Jul 2003</u> . . Evercat (fmt)
Animated Traces [Ben Fry]



<u>http://benfry.com/traces/</u>

FR Français (France) 🙎 📋 🖶 Java - WelcomePageDispatchAction, java - Eclipse SDK File Edit Navigate Search Project Run Window Help 📬 • 🔚 🖻 🗄 🧶 • 🍫 • 🔘 • 💁 • 🗄 🍄 🞯 • 🗄 🥭 🖋 • 🗄 🧏 - 🎘 - 🏷 -😭 🐉 Java 指 Hierarchy) 🕞 🔄 🎽 🗖 🖬 🚰 Fragment Comparison 🛛 🚺 LoginDispatchAction. × 増 Package Explorer 🖾 CollectInformationDi 🖃 对 Struts demo Previous Version Current Version 🕮 src // End of user imports // End of user imports ~ 😑 🥵 WEB-INF/src | 🗄 🖶 com.mia_software.booster.page_flow_example public class WelcomePageDispatchAction public class WelcomePageDispatchAct 🗄 🖶 com.mia_software.booster.page_flow_example 😑 🖶 com.mia_software.booster.page_flow_example // associated forward definitions // associated forward definitio 🗄 🖳 CollectInformationDispatchAction.java public final static String COLLECT public final static String COLL 🗄 🖳 CollectInformationForm.iava public final static String LOGOUTAN public final static String LOGO 🗄 – 🛄 ResultDispatchAction, java public final static String TEST1 TC public final static String TEST 🗄 🛄 ResultForm.java 🖮 🛄 Test1DispatchAction. java // inherited forward definitions // inherited forward definition 🖻 🖳 Test1Form.java 🗄 – 🛄 WelcomePageDispatchAction, java 🗄 🖳 WelcomePageForm.java // dispatch action methods declarat // dispatch action methods decl 🗄 🔚 com.mia_software.struts.generic.back public ActionForward enter(ActionMa public ActionForward enter(Acti 😟 🖶 com.mia_software.struts.generic.form ActionForward actionForward = r ActionForward actionForward < > if (form != null) { if (form != null) { WelcomePageForm current WelcomePageForm currentForn 🙋 Generation Results 🛛 N // execute code on exit of // execute code on exit 🖃 🚯 Results (30) currentForm.onExit(); currentForm.onExit(); 🗄 🐨 🔜 ResultDispatchAction, iava CollectInformationForm int i ENTERXXX=0; 🗄 🚽 🔂 Test1DispatchAction.java CollectInformationForm coll 😑 🔛 WelcomePageDispatchAction.java // Start of user code : 🗟 Generated Code // Start of user code : ret Dianual Code 🗟 Generated Code Manu 🗸 Sort 🐻 Gener 🗸 Group by Status 💮 Manu 💿 File Name 🗟 Genei Relative Path 🗋 Manu Full Path 💩 Genei

•

💮 CollectInf

Compare with Previous Version

v

<

Show In

>

<

>

Diff style: Side-by-side 🗸

/home/toddw/src/sshconsole-read-only/content/sshconsole.js

Files Changed:

1. sshconsole.js: 1 change [1]

E0 Ener Midden (Energy)									
	term - nr.: \/T100/00_24_"term"\.	ndden E4	terr - pri //TIOO/RO 24 //terr//):						
51	_term = new viido(ob, 24, term);	51	_term = new V100(00, 24, "term");						
52	//_term.debug_ = 1;	52	//_term.debug_ = 1;						
53	_term.curs_set(true, true, _term_box_etement);	53	_term.curs_set(true, true, _term_box_etement);						
94 EE		54	_term.noecho();						
50	// Deplace the selected function with our sum, this is called	55	// Deplace the concertsh function with our own, this is called						
50	// Replace the go_getch_ function with our own, this is catted	50	// Reptate the go_gettin_ function with our own, this is tatted						
51	// for every keypress that is passed through the terminat to the	57	// for every keypress that is passed through the terminat to the						
50	// remote server. The character is atready converted into the	50	// remote server. The character is acready converted into the						
59	// required viido character sequence(s).	59	// required vitoo character sequence(s).						
60	VIIOU.go_getch_ = Tunction() {	60	Vilde.go_getch_ = Tunction() {						
61	var vt = vi100.the_vt_;	61	var vt = v1100.the_vt_;						
62	1T (VT === underined) {	62	1T (VT === somevalue) {						
63	return;	63	return;						
64		64	}						
65	<pre>var ch = vt.key_butshift();</pre>	65	<pre>var ch = vt.key_butshift();</pre>						
66	//dump("go_getch_:: ch: '" + ch + "'\n");								
67	if (ch === undefined) {	66	if (ch === undefined) {						
68	return;	67	return;						
69	}	68	}						
70	if (vt.echo_ && ch.length == 1) {	69	if (vt.echo_ & ch.length == 1) {						
71	vt.addch(ch);	70	vt.addch(ch);						
		71	vt.refres();						
72	}	72	}						
73	if (_ssh_channel) {	73	if (_ssh_channel) {						
74	_ssh_channel.sendStdin(ch);	74	_ssh_channel.sendStdin(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_textbox");</pre>						
79	var connectionText;	79	var connectionText;						
80	<pre>if ('connectionText' in window.arguments[0]) {</pre>	80	<pre>if ('connectionText' in window.arguments[0]) {</pre>						
81	<pre>connectionText = window.arguments[0].connectionText;</pre>	81	<pre>connectionText = window.arguments[0].connectionText;</pre>						
82	} else {	82	} else {						
	174 lines hidden [Expand]								

^

History Flow (Viégas et al)



"Abortion"

authors

Zundark

The Epopt

B4hand

Shubenstein

Maverie149

Theanthrope

Drehmword

Comembert

Hephoestos

MyRedDice

Kingturtie

2001

from Wikipedia

COLOR 💒 group 📜 individual 🛛 🔛 text changes 🔉 text age SPACING O date O versions



Abortion

(Revision as of 22:56 4 Jun 2003)

"Abortion," in its most commonly used se refers to the deliberate early termination pregnancy, resulting in the death of the gr fetus, [1] Medically, the term also refers t early termination of a pregnancy by natur ("spontaneous abortion" or miscarriage, w 1 in 5 of all pregnancies, usually within the weeks) or to the cessation of normal grow body part or organ. What follows is a disci the issues related to deliberate or "induceabortion.

Methods

Depending on the stage of pregnancy an a performed by a number of different metho a chemical abortion is the usual method, t mifepristone is usually the only legal meth although research has uncovered similar e from methotrexate and misoprostal. Conc with chemical abortion and extending up u around the fifteenth week suction-aspiration vacuum abortion is the most common app replacing the more risky dilation and curet C). From the fifteenth week up until aroun eighteenth week a surgical dilation and ex (D & E) is used.

As the fetus size increases other technique be used to secure abortion in the third trip premature expulsion of the fetus can be in with prostaglandin, this can be coupled with injecting the amniotic fluid with saline or u solution. Very late abortions can be broug by the controversal intact dilation and extr & X) or a hysterotomy abortion, similar to caesarian section.

The controversy

The morality and legality of abortion is a l important topic in applied ethics and is als discussed by legal scholars and religious p Important facts about abortion are also re by sociologists and historians.

Abortion has been common in most societ although it has often been opposed by sor institutionalized religions and governments century politics in the <u>United States</u> and <u>En</u> abortion became commonly accepted by the 20th century. Additionally, abortion is accepted in China. India and other populo countries. The Catholic Church remains o the procedure, however, and in other cour notably the United States and the (predom Catholic) Republic of Ireland, the controve extremely active, to the extent that even t of the respective positions are subject to I debate. While those on both sides of the are generally peaceful, if heated, in their i of their positions, the debate is sometimes characterized by violence. Though true of sides, this is more marked on the side of t opposed to abortion, because of what they the gravity and urgency of their views.

The central question

2003

The central question in the abortion debat clash of presumed or perceived rights. On hand, is a fetus (sometimes called the "un pro-life/anti-abortion advocates) a human with a right to life, and if so, at what point pregnancy does the fetus become human? other hand, is a fetus part of a woman's b



Visualizing Conversation

Many dimensions to consider:

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

Interesting cross-products:

- What x When \rightarrow Topic "Zeitgeist"
- Who x Who \rightarrow Social network
- Who x Who x What x When \rightarrow Information flow

Naming Names

Names used by major presidential candidates in the series of Democratic and Republican debates leading up to the Iowa caucuses.



Usenet Visualization (Viégas & Smith)

Show correspondence patterns in text forums Initiate vs. reply; size and duration of discussion



back to newsgroups

Week of Oct 21, 2001

......

0.....

8

subject, subject	# of posts
ednesday Spooker ASF	21
ET#3 Anyone for breakfast	
anny Side Up ASF ()	
aturday Ensemble and WET	
h noi Watch outli ASF	
iursday Combo-Post WET #	
ie Yellow Rose InnA gitt to	
ET #1 JBP The First Time	
e Love the Earth ASF	15
onday Spocker "The Sight"	16
	14
teberge "Le Vent Se Leve"	14
oliday Tog #3 🜖	13
pooker da Jour 🦙 🤅	13
eginning ASF Short and	
econd Try A Kalle for Suzy	12
ome On a Safari With Me	11
Jesday Spooker ASF	11
urses, Foiled Again ASF	
alloween Togs Take Two	
eauty of the Fury Jim Warren	
hought i saw	
ednesday Evening at the Con	
econd Try A Kalle for Suzy	
ank Was A Monster ASF	

-

subject	# of posts	
Sunday Twofer ASF)	9	
Chopsticks/A Jilly fake		
Oh not Trouble in Discworld!	7	
WETyourthirst1_ASF	6	
A pretty for youReposted fro		
Saturday Spooker ASF		
Sample Previous install Upgr	4	
Tennessee weather tonite	4	
WET - Well I am not smiling!	4	
Somethin' mushy <ast></ast>		
Getting seasonal with workin		
A Haunted House)		
do you wonder what debits be		
Question Ethics of posters in	3	
For Jerry		
Olu's Tribe - slightly rated		
WET - Glass Bottles	3 1	
Peace Train <asf></asf>	2	
Arrival at Stewart Island II	2	
WET 195 Wrap-up	2	
Cat O'Lantern	2	
Put a Spell on You (Happy H	2	
Goodbye to Summer - A Timel	2	
Two Pumpkins In A Strange B	2	
Still Heading South II	2	
WET- Frank Sinatra - The Man.	2	
WET Autumn	2	
Purple Martin ASF	2	
Opposites Attract	2	
Time	2	
	14	- M

1

author:

jillyb@mail.com



1.1

Mountain (Viégas)

Conversation by person over time (who x when).



Themail (Viégas et al)



One person over time, TF.IDF weighted terms

Enron E-Mail Corpus

[Heer]

👙 Enron Corpus Viewer







👙 Enron Corpus Viewer





+ about + permalink	+ SELECT ALL COUNT	RIES AUSTRALIA	AUSTRIA	CANADA	FRANCE	DEUTSCHLAND	INDIA	Italia	NEW ZEALAND	ESPANA	U.K .		<mark>≪</mark> U.S.
Pakistan's Clinton criticises			n ses	Cricket ref: Pakistan police fled during			Was Fallon funny?	CONCERT REVIEW: Britney Yet to Hit Stride	Spotlight still follows the Jonas Brothers	The R To Ba Rush	ush sh	The Strengths and Weaknesses of the Chandra Levy Case	
international demoliting arena demolitional E.Jerusa		litio	How Could Rihanna Take Back Chris Brown?				Jason & Molly: 'Fighting Like Hell to Make It Work'	'Watchmen' tickets are hot: Have you bought yours yet?	Marines: Multiple errors caused	Celebrating victory, Villaraigosa acknowledge challenges af	Chester Stiles found guilty on all counts		
		E.Jeru	erusal ambus		n		Economy to Dominate Annual Chinese Gathering	Cops 'High Called to Musi Octo gets -Home 8 cast Times: Report Mile	School Catin the Hat stresses a new of reading to Boone students y Parts Wrep to Adapt and Close warrs, Benamed	San Diego crash Obama: Contractin	Fed spending bill contains billions in earmarks for L	Supreme Court Divided Over Judicial Bias Case Tapped to Head FCC	
			em		a .		Cirordi	Michael Jackson seeks comeback	U2 Talk "Horizon" Follow Up, Spider-Man Musical in 'Arerica	IC investiget Meganicrea /S divers frearing	g overhaul to save	Voter (s) elect Emanuel's likely replacement —	Autorit Budget parier Batege Contrare in Berger Semate race is Provide race
British	Command er says Iran missiles can reach Israel atom sites	World court issues arrest warrant for Sudan's Bashir	Afghan Bomb Kills Election Three Commission Canadian Says Early Servicemen Vote Not in Southern	Bomb Kills	Coast	Man Says hip		in London	Rolling Stone Read Li Bounds, F Burley, Machigen	disia in the second sec	\$40 billion a year	no one cares	now an online and the second s
Prime Minister to				Guard Calls Off	ny, LA bothere Close year		Oil Gains a Second Day of Speculation	Chrysler on spends \$55 per car on l	5566 Apple 5566 Introdu Mac de s, compu	uces new TS esktop an iters thi	el and MC: What e they inking?	Senators Baucus Sets push to Dehvery of boost FDA Healthcare food safety system	
Address US Congress		Has Pakistan	POSSIDIE South Mexice Africa: troops United drug v States May drug v	Afgnanistan o Gandhi kin enter var memorabil	Search for Missing	Suns- Heat Preview	Beckham's crafty new oan deal exposes the sign Kurt	Stimulus	Canada	In-depth review: Kindle 2, the Apple	CEBIT-IT industry will be back, Schwarzen	Phew! Asteroid's passing was a	AIDS affecting Formar HP oxider adults leader disgnosed with cancer NG doctor Concerted being Effect Mandud
Interim Leader Takes Over in	China, US opposed to	, US become become Boycotus may cit the central front? Conference control and the central front? Conference control control and control an	Boycott UN Racism Conference Smaller del boost amid economy w	ia Man hoping for 'world's ^{ses} best job'	Players, Fellow Boater	truth Warner? Pacers- Trail Blazers Preview Golde n State Free-sent Bit Abulates/ State State State Red Wings- Avalanche Preview LA Clipp ers Des to State Market State	Bernanke's H AlG Blast May B Mean More Curbs on	ollywood mulls fe without lockbuster BS Gets	states and the second s	Mobile- phone discover Competition orbiting to Point Schut	cosmic near-miss ed Microsoft is informally backing a new search seque same at live	Track online track online PEA Paralls Canadar Wear the claims linked to one tooth	
Guinea-Bissau	NKorean missile launch: envoy		Germany: 2 still missing after Bad building altimeter factor in Netherla plane cra	Cameroon: Govt, Church to Sitare Cost of Repeter State tor in meter a tor in hertands ne crash			A Clippo A Considers Statis	Stars Risk, Concentration Ukraine May Miss Deadline for Gazprom Gas Payment	No orise interview intervi	PC Shipments Expected To Suffer Sheet PC Shipments Expected To Suffer Sheet Party like	9 Sales Slow 9 The Foundry Company Becomes "Globalfour dries". Posteria Aguin	n Article of Article A	MY country Substantian of the second
Wednesday March 4, 2009 8:09				1500		ELECT ALL CATEGORIES	DUCIUERO	TECHNOLOGY	coopte		UARIFIED STANDARD		
	YEST. TODAY NOW					LESS	Than 10 minutes ago		BUSINESS	TECHNOLUGY	SPORIS	ENTERTHIN	
12:00						м	ore than 1 hour ago						

NewsMap: Google News Treemap (Marcos Weskamp)

Tips: Document Contents

Understand your task, and handle high dimensionality accordingly...

Visually: Word position, browsing, brushing + linking Semantically: Word sequence, hierarchy, clustering Both: Spatial layout reflect semantic relationships

Role of Interaction:

Sufficient language model to enable visual analysis cycles Allow modifications to the model: custom patterns for expressing contextual or domain knowledge