

Text and Document Visualization 2



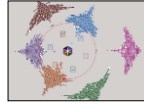
CS 7450 - Information Visualization
November 2, 2016
John Stasko

Learning Objectives

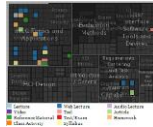


- List different queries/tasks often needed on document collections
- List various analytic metrics often calculated on documents
- List different aspects of documents often visualized
- Explain vector space document analysis (similarity calculation, search)
- Explain TFIDF
- Describe visual representation used by and contributions of these systems
 - Themail, PaperLens, Jigsaw, ThemeScape/IN-SPIRE, ThemeRiver

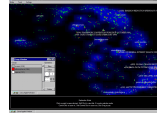
This Week's Agenda



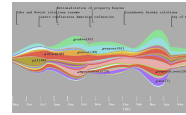
Visualization for IR
Helping search



Visualizing text
Showing words,
phrases, and
sentences



Visualizing document sets
Words & sentences
Analysis metrics
Concepts & themes



Last Time

Fall 2016

CS 7450

3

Today's Agenda



- Move to collections of documents
 - Still do words, phrases, sentences
 - Add
 - More context of documents
 - Document analysis metrics
 - Document meta-data
 - Document entities
 - Connections between documents
 - Documents concepts and themes

Fall 2016

CS 7450

4

Related Topic - Sensemaking



- Sensemaking
 - Gaining a better understanding of the facts at hand in order to take some next steps
 - (Better definitions in VA lecture)
- InfoVis can help make a large document collection more understandable more rapidly

Your HW



- What tasks/goals/questions could be expected with the Amazon reviews data set?
- Let's generate a list...

Questions/Tasks



How many reviewers recommend the TV?
How do the reviewers rate the different features of the TV?
How long has a reviewer had the TV?
How many five star reviews are there for the TV?
What are the favorite and least favorite factors for the TV?
Was there anything associated with the word 'annoying'?
How long did it take to get a negative review?
Do they recommend it to others?
How does it compare to other TVs?
What do users think about the size?
Is it good value for money?
For bad reviews, what are the frequent negative reviews?
Do users have a particular use for the TV?
If it has problems, are they easy to fix?
Was delivery quick and easy?
How many reviews speak about a specific feature?
How was the setup process?
How many reviews talk about the purchase and service experience?
What is the average rating?
Is it a good TV or not?
How is the feature X rated?
Have the review sentiments changed over time?
What do they think of the brand / compare it to others?

Fall 2016

CS 7450

7

Overlaps & Similarities



- Are some items in our list in the same "category"?
 - Can we generalize a little and narrow the list down to some core questions/tasks?

Fall 2016

CS 7450

8

Questions/Tasks



Is this a good TV? (summarize ratings)
What are the positive/negative features?
Characterize product features vs other stuff.
What do people think about X?
Comparison to other TVs.
What are the defects/problems?
What are the views over time?

Fall 2016

CS 7450

9

Evaluate a Vis



- Use that list to evaluate a visualization for this problem

Fall 2016

CS 7450

10

Example Tasks & Goals



- Which documents contain text on topic XYZ?
- Which documents are of interest to me?
- Are there other documents that are similar to this one (so they are worthwhile)?
- How are different words used in a document or a document collection?
- What are the main themes and ideas in a document or a collection?
- Which documents have an angry tone?
- How are certain words or themes distributed through a document?
- Identify "hidden" messages or stories in this document collection.
- How does one set of documents differ from another set?
- Quickly gain an understanding of a document or collection in order to subsequently do XYZ.
- Understand the history of changes in a document.
- Find connections between documents.

Various Document Metrics



- Different variables for literary analysis
 - Average word length
 - Syllables per word
 - Average sentence length
 - Percentage of nouns, verbs, adjectives
 - Frequencies of specific words
 - Hapax Legomena – number of words that occur once

Vis

Each block represents a contiguous set of words, eg, 10,000 words

Do partial overlap in blocks for a smoother appearance

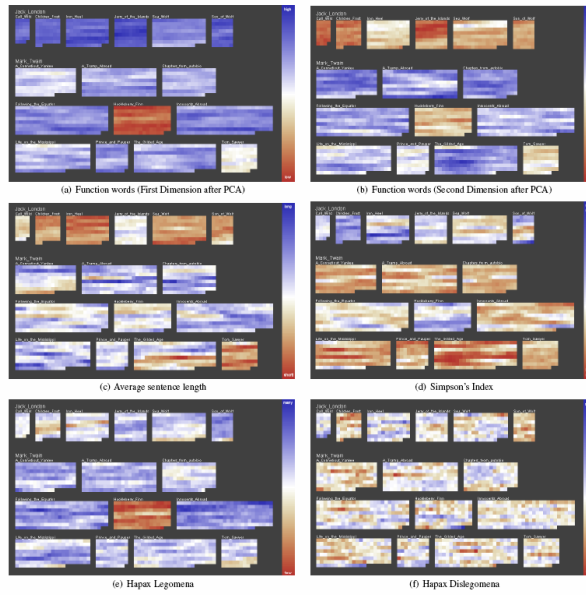


Figure 2: Fingerprints of books of Mark Twain and Jack London. Different measures for authorship attribution are tested. If a measure is able to discriminate between the two authors, the visualizations of the books that are written by the same author will equal each other more than the visualizations of books written by different authors. It can easily be seen that this is not true for every measure (e.g. Hapax Dislegomena). Furthermore, it is interesting to observe that the book *Huckleberry Finn* sticks out in a number of measures as if it is not written by Mark Twain.

The Bible

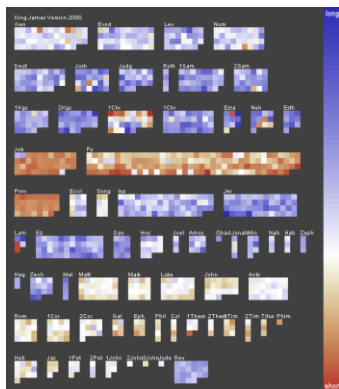


Figure 4: Visual Fingerprint of the Bible. Each pixel represents one chapter of the bible and color is mapped to the average verse length. Interesting characteristics such as the generally shorter verses of the poetry books, the inhomogeneity of the 1. Book of Chronicles or the difference between the Old Testament and the New Testament can be perceived.

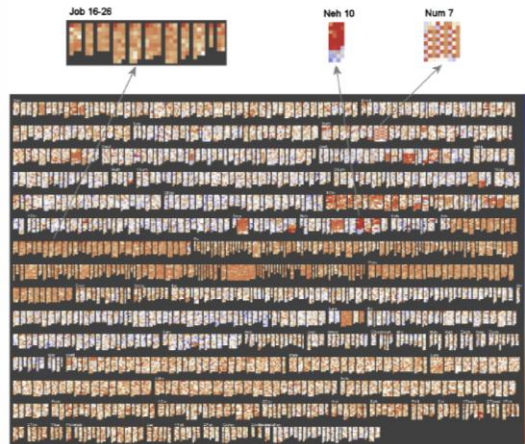


Figure 5: Visual Fingerprint of the Bible. More detailed view on the bible in which each pixel represents a single verse and verses are grouped to chapters. Color is again mapped to verse length. The detailed view reveals some interesting patterns that are camouflaged in the averaged version of fig. 4.

Follow-On Work



- Focus on readability metrics of documents
- Multiple measures of readability
 - Provide quantitative measures
- Features used:
 - Word length
 - Vocabulary complexity
 - Nominal forms
 - Sentence length
 - Sentence structure complexity

Oelke & Keim
VAST '10

Visualization & Metrics



		Voc. Difficulty	Word Length	Nominal Forms	Sent. Length	Compl. Sent. Struc.
(a)	The intention of TileBars [9] is to provide a compact but yet meaningful representation of Information Retrieval results, whereas the FeatureLens technique, presented in [5], was designed to explore interesting text patterns which are suggested by the system, find meaningful co-occurrences of them, and identify their temporal evolution.					
(b)	This includes aspects like ensuring contextual coherency, avoiding unknown vocabulary and difficult grammatical structures.					

Figure 5: Two example sentences whose overall readability score is about the same. The detail view reveals the different reasons why the sentences are difficult to read.

Uses heatmap style vis (blue-readable, red-unreadable)

Interface

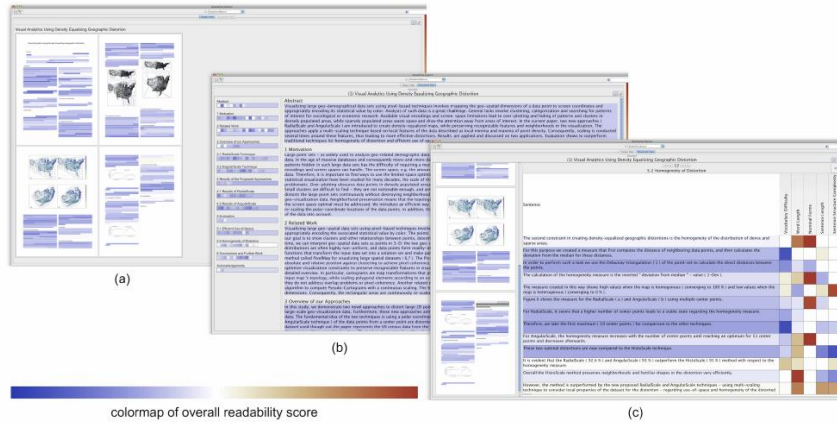


Figure 3: Screenshot of the VisRA tool on 3 different aggregation levels. (a) Corpus View (b) Block View (c) Detail View. To display single features, the colormap is generated as described in section 3.4 and figure 2.

Fall 2016

CS 7450

17

Their Paper (Before & After)

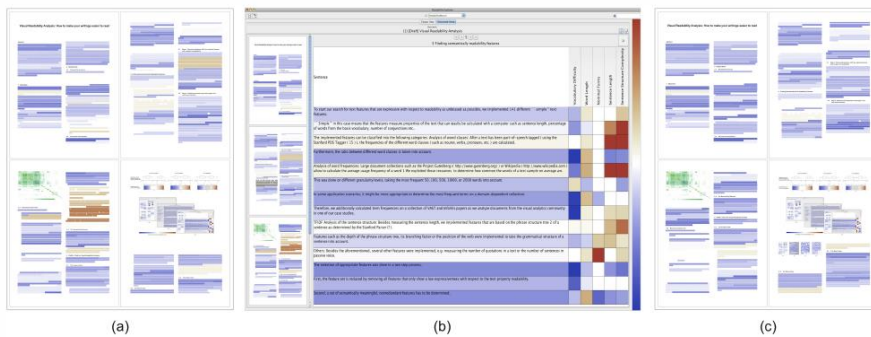


Figure 6: Revision of our own paper. (a) The first four pages of the paper as structure thumbnails before the revision. (b) Detail view for one of the sections. (c) Structure thumbnails of the same pages after the revision.

Fall 2016

CS 7450

18

Comment from the Talk



- In academic papers, you want your abstract to be really readable
- Would be cool to compare rejected papers to accepted papers

Fall 2016

CS 7450

19

Bohemian Bookshelf

Video



Serendipitous browsing



Fall 2016

CS 7450

Thudt et al
CHI '12

20

Themail



- Visualize one's email history
 - With whom and when has a person corresponded
 - What words were used
- Answer questions like:
 - What sorts of things do I (the owner of the archive) talk about with each of my email contacts?
 - How do my email conversations with one person differ from those with other people?

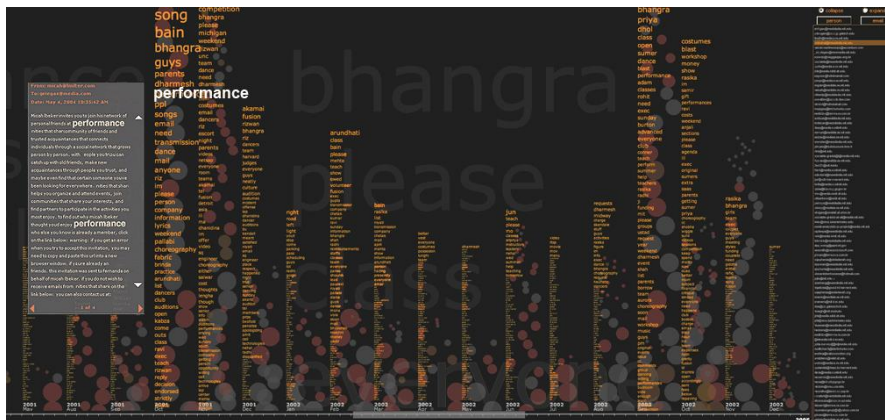
Viégas, Golder & Donath
CHI '06

Fall 2016

CS 7450

21

Interface



Fall 2016

CS 7450

22

Characteristics



- Text analysis to seed visualization
- Monthly & yearly words



Figure 2: Expanded view of Themail showing the sporadic nature of a relationship. “Blank” spaces between columns of words stand for months when no messages were exchanged between the user and the selected email contact.

Fall 2016

CS 7450

23

Query UI

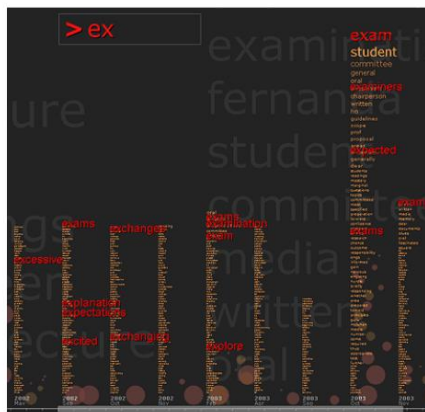


Figure 3: Searching for words in Themail. Here the user has typed “ex” (at the top of the screen) and Themail has highlighted (in red) all the monthly words starting with these characters.

Fall 2016

CS 7450

24

PaperLens



- Focus on academic papers
- Visualize doc metadata such as author, keywords, date, ...
- Multiple tightly-coupled views
- Analytics questions
- Effective in answering questions regarding:
 - Patterns such as frequency of authors and papers cited
 - Themes
 - Trends such as number of papers published in a topic area over time
 - Correlations between authors, topics and citations

Fall 2016

CS 7450

Lee et al
CHI '05 Short

25

PaperLens

Video



- a) Popularity of topic
- b) Selected authors
- c) Author list
- d) Degrees of separation of links
- e) Paper list
- f) Year-by-year top ten cited papers/ authors – can be sorted by topic

Fall 2016

CS 7450

26

More Document Info



- Highlight entities within documents
 - People, places, organizations
- Document summaries
- Document similarity and clustering
- Document sentiment

Fall 2016

CS 7450

27

Jigsaw



- Targeting sense-making scenarios
- Variety of visualizations ranging from word-specific, to entity connections, to document clusters
- Primary focus is on entity-document and entity-entity connection
- Search capability coupled with interactive exploration

Stasko, Görg, & Liu
Information Visualization '08
Görg et al
IEEE TVCG '13

Fall 2016

CS 7450

28

Document View



Wordcloud overview

Doc List

Document summary

Selected document's text with entities identified

analysis analysts **analytic** animation based **cognition** design discuss display **evaluation**

framework **information** infovis **interaction** level localization paper research

systems tasks techniques video visual **visualization** visualizations

Documents

- 1 infovis00-895991
- 0 infovis01-963277
- 0 infovis03-1249027
- 1 infovis04-1382902
- 1 infovis05-1532136
- 1 infovis07-4376134
- 0 infovis07-4376144
- 2 infovis08-4658127
- 0 infovis08-4658129
- 2 infovis08-4658146
- 0 infovis09-5290708
- 0 infovis09-5290708
- 0 infovis09-5290708
- 0 vast07-4389006
- 0 vast07-4389013
- 0 vast09-5332596
- 1 vast09-5333878

Summary: Evaluating visual analytics systems for investigative analysis. Deriving design principles from a case study Despite the growing number of systems providing visual analytic support for investigative analysis, few empirical studies of the potential benefits of such systems have been conducted, particularly controlled, comparative evaluations. Determining how such systems foster insight and sensemaking is important for their continued growth and study, however. Furthermore, studies that identify how people use such systems and why they benefit (or not) can help inform the design of new systems in this area. We conducted an evaluation of the visual analytics system Jigsaw employed in a small investigative sensemaking exercise, and we compared its use to three other more traditional methods of analysis. Sixteen participants performed a simulated intelligence analysis task under one of the four conditions. Experimental results suggest that Jigsaw assisted participants to analyze the data and identify an embedded threat. We describe different analysis strategies used by study participants and how computational support (or the lack thereof) influenced the strategies. We then illustrate several

Fall 2016

CS 7450

29

List View

Entities listed by type



Concept

- interaction
- evaluation
- insight
- visual analytics
- case study
- cognition
- color
- navigation
- animation
- category
- document
- dynamic query
- filter
- focus+context
- hierarchy
- intelligence analysis
- metrics
- perception
- social
- software visualization
- text
- theory
- time series
- treemap
- awareness
- biometrics
- brushing
- business
- cluster
- collaboration
- data mining
- database
- education
- elbow
- geographic
- graph
- hardware
- high-dimensional data

author

- Spence, B.
- Springue, D.W.
- Sorenson, T.C.
- Sprong, N.
- Stadler, P.F.
- Stasko, J.
- Steed, C.A.
- Stien, C.
- Stocking, K.
- Stofel, A.
- Stoth, C.
- Storey, M.A.D.
- Strasser, T.
- Strayer, D.
- Stroholtz, H.
- Strohmann, P.J.
- Stuckey, P.
- Stulen, F.
- Sturtebeck, E.P.
- Sturtz, D.
- Sui, H.
- Sudhanto, A.
- Suh, B.
- Sullivan, T.
- Suma, E.
- Summers, K.L.
- Sunm, J.
- Swan, J.E.
- Swindle, C.
- Szardis, N.
- Takedama, Y.
- Tal, A.
- Talbot, J.
- Tan, D.S.
- Tan, R.
- Tanase, T.
- Tandon, S.
- Tang, D.
- Tanm, E.
- Tatui, A.
- Tavanti, M.

year

- 1995
- 1996
- 1997
- 1998
- 1999
- 2000
- 2001
- 2002
- 2003
- 2004
- 2005
- 2006
- 2007
- 2008
- 2009

conference

- Infovis
- VAST

Fall 2016

CS 7450

30

Document Cluster View

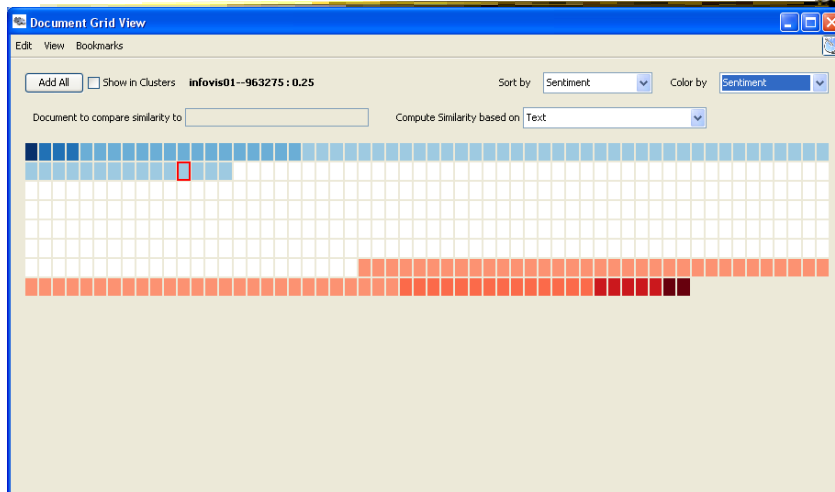


Fall 2016

CS 7450

31

Document Grid View



Here showing sentiment analysis of docs

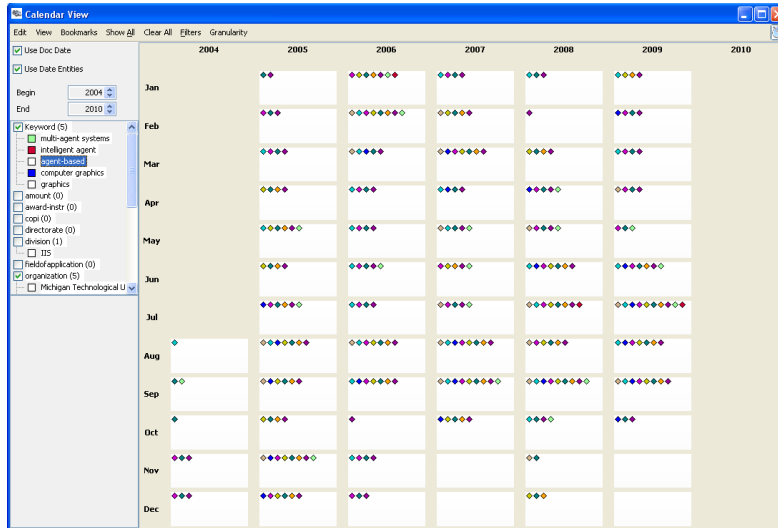
Fall 2016

CS 7450

32

Calendar View

Temporal context
of entities & docs

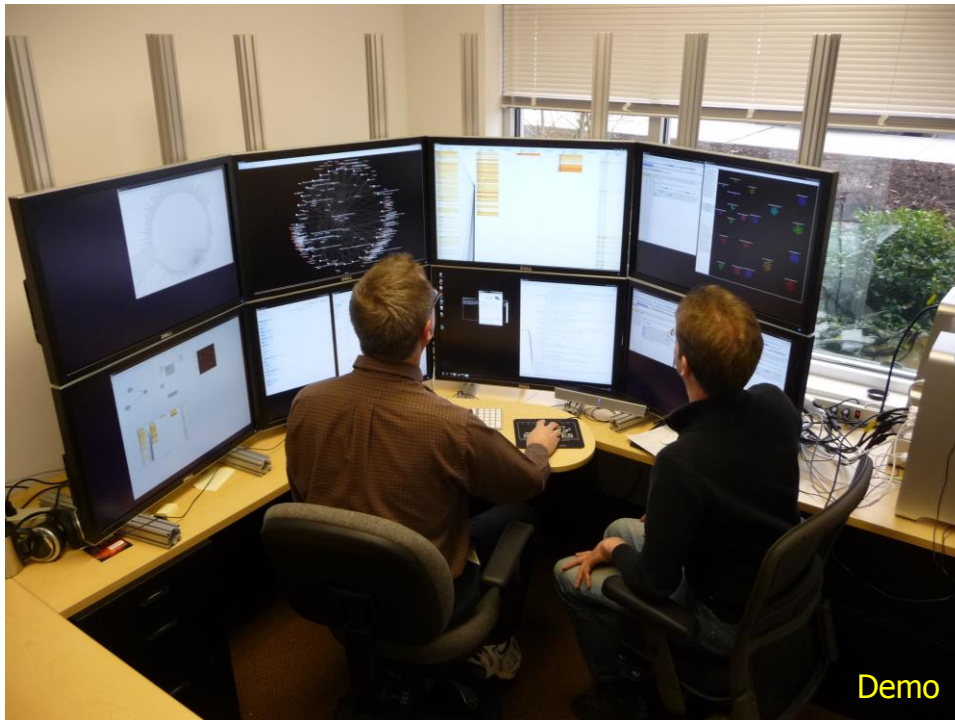


Fall 2016

CS 7450

Video

33



Jigsaw



- More to come on Visual Analytics day...

Fall 2016

CS 7450

35

FacetAtlas



- Show entities and concepts and how they connect in a document collection
- Visualizes both local and global patterns
- Shows
 - Entities
 - Facets – classes of entities
 - Relations – connections between entities
 - Clusters – groups of similar entities in a facet

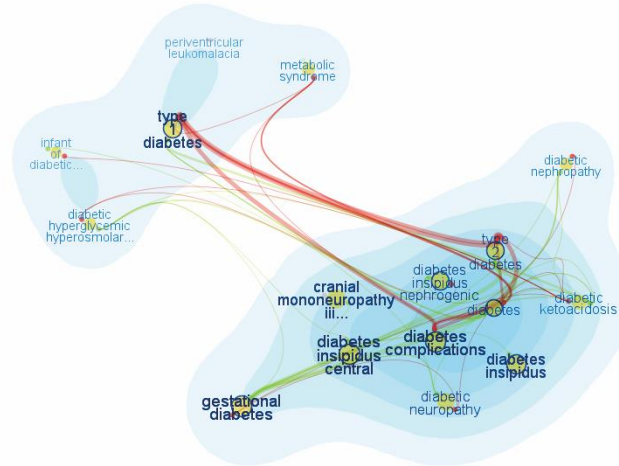
Cao et al
TVCG (InfoVis) '10

Fall 2016

CS 7450

36

Visualization



Fall 2016

CS 7450

37

Up to Higher Level



- How do we present the contents, semantics, themes, etc of the documents
 - Someone may not have time to read them all
 - Someone just wants to understand them
- Who cares?
 - Researchers, fraud investigators, CIA, news reporters

Fall 2016

CS 7450

38

Vector Space Analysis



- How does one compare the similarity of two documents?
- One model
 - Make list of each unique word in document
 - Throw out common words (a, an, the, ...)
 - Make different forms the same (bake, bakes, baked)
 - Store count of how many times each word appeared
 - Alphabetize, make into a vector

Fall 2016

CS 7450

39

Vector Space Analysis



- Model (continued)
 - Want to see how closely two vectors go in same direction, inner product
 - Can get similarity of each document to every other one
 - Use a mass-spring layout algorithm to position representations of each document
- Some similarities to how search engines work

Fall 2016

CS 7450

40

Wiggle



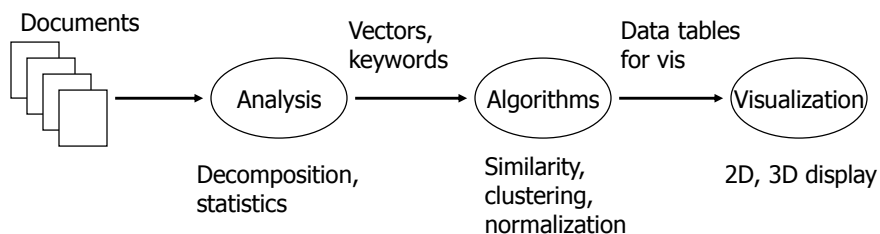
- Not all terms or words are equally useful
- Often apply TFIDF
 - Term frequency, inverse document frequency
- Weight of a word goes up if it appears often in a document, but not often in the collection

Fall 2016

CS 7450

41

Process



Fall 2016

CS 7450

42

VIBE System



- Smaller sets of documents than whole library
- Example: Set of 100 documents retrieved from a web search
- Idea is to understand contents of documents relate to each other

Olsen et al
Info Process & Mgmt '93

Fall 2016

CS 7450

43

Focus



- Points of Interest
 - Terms or keywords that are of interest to user
 - Example: cooking, pies, apples
- Want to visualize a document collection where each document's relation to points of interest is show
- Also visualize how documents are similar or different

Fall 2016

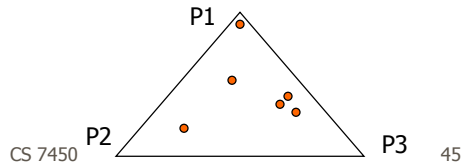
CS 7450

44

Technique



- Represent points of interest as vertices on convex polygon
- Documents are small points inside the polygon
- How close a point is to a vertex represents how strong that term is within the document

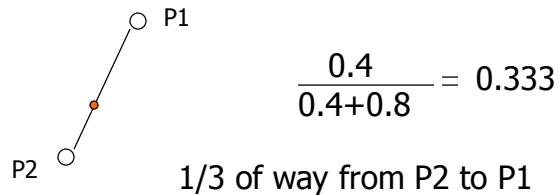


Fall 2016

Algorithm



- Example: 3 POIs
- Document (P1, P2, P3) (0.4, 0.8, 0.2)
- Take first two



Fall 2016

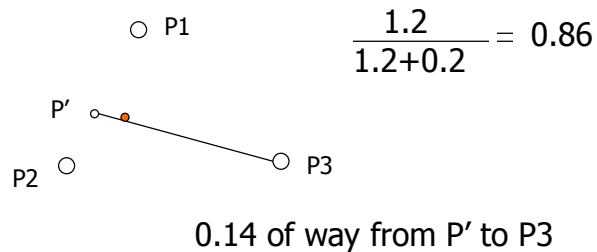
CS 7450

46

Algorithm



- Combine weight of first two 1.2 and make a new point, P'
- Do same thing for third point

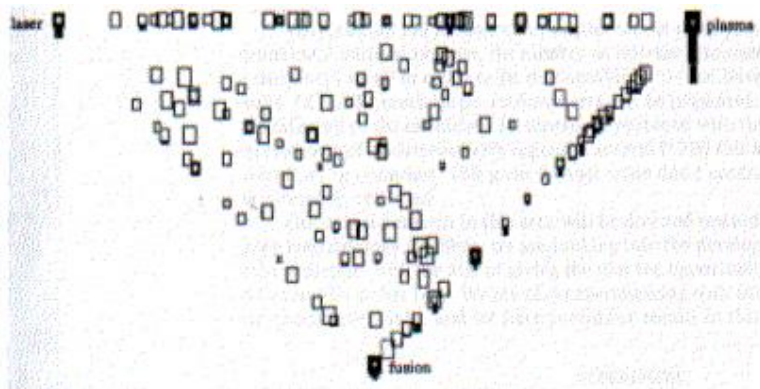


Fall 2016

CS 7450

47

Sample Visualization



Fall 2016

CS 7450

48

VIBE Pro's and Con's



- Effectively communications relationships
- Straightforward methodology and vis are easy to follow
- Can show relatively large collections
- Not showing much about a document
- Single items lose “detail” in the presentation
- Starts to break down with large number of terms

Visualizing Documents



- VIBE presented documents with respect to a finite number of special terms
- How about generalizing this?
 - Show large set of documents
 - Any important terms within the set become key landmarks
 - Not restricted to convex polygon idea

Basic Idea



- Break each document into its words
- Two documents are “similar” if they share many words
- Use mass-spring graph-like algorithm for clustering similar documents together and dissimilar documents far apart

Fall 2016

CS 7450

51

Kohonen’s Feature Maps



- AKA Self-Organizing Maps
- Expresses complex, non-linear relationships between high dimensional data items into simple geometric relationships on a 2-d display
- Uses neural network techniques

Lin
Visualization '92

Fall 2016

CS 7450

52

Map Attributes



- Different, colored areas correspond to different concepts in collection
- Size of area corresponds to its relative importance in set
- Neighboring regions indicate commonalities in concepts
- Dots in regions can represent documents

Fall 2016

CS 7450

53

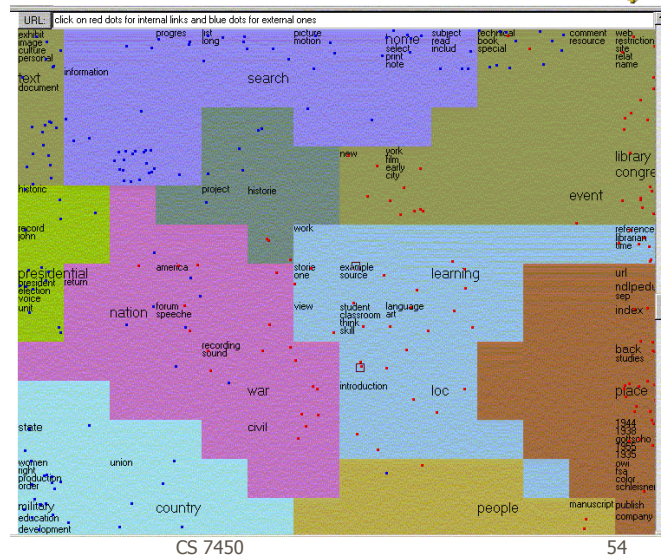
More Maps



Interactive demos

Xia Lin

Fall 2016



CS 7450

54

Work at PNNL



- Group has developed a number of visualization techniques for document collections
 - Galaxies
 - Themescapes
 - ThemeRiver
 - ...

Wise et al
InfoVis '95

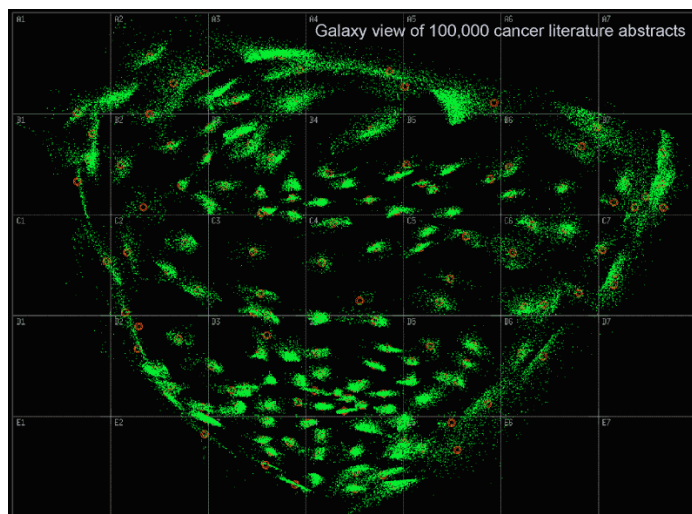
Fall 2016

CS 7450

55

Galaxies

Presentation of documents where similar ones cluster together



Fall 2016

CS 7450

56

Themescapes



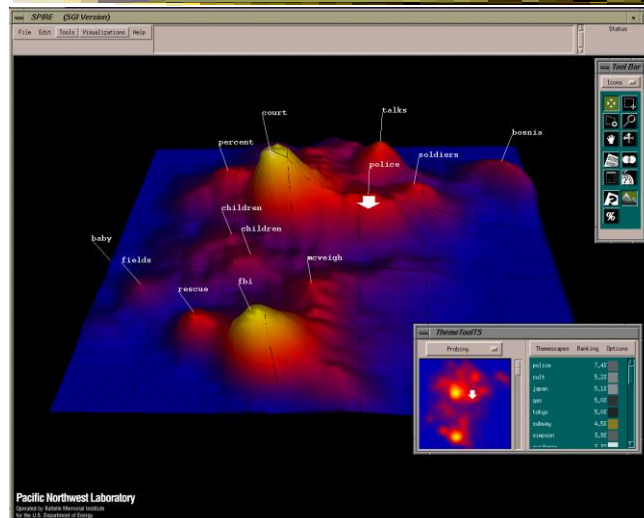
- Self-organizing maps didn't reflect density of regions all that well -- Can we improve?
- Use 3D representation, and have height represent density or number of documents in region

Fall 2016

CS 7450

57

Themescape



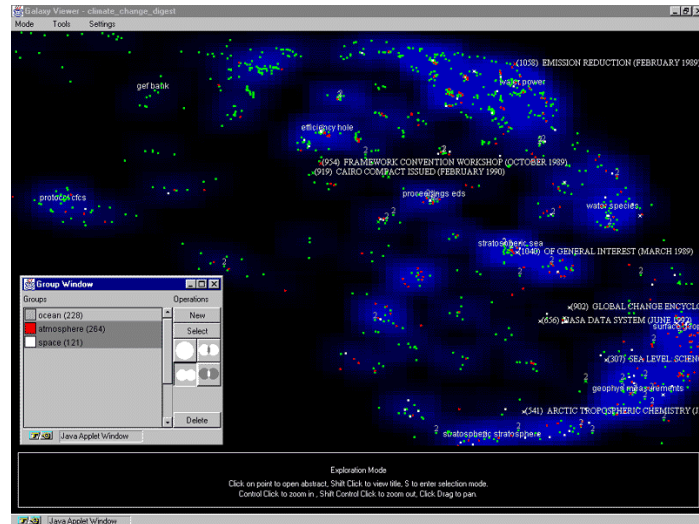
Video

Fall 2016

CS 7450

58

WebTheme



Fall 2016

CS 7450

59

Related Topic



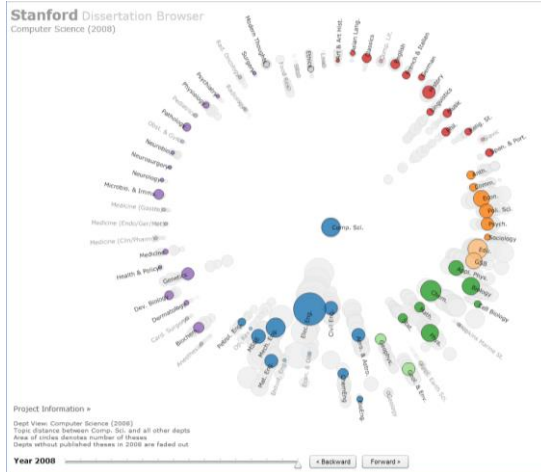
- Maps of Science
- Visualize the relationships of areas of science, emerging research disciplines, the impact of particular researchers or institutions, etc.
- Often use documents as the “input data”

Fall 2016

CS 7450

60

Stanford Diss. Browser



9,000 Stanford PhD theses

Rather than overall semantic map, you choose a focus and all update to show their relationship

Demo at

<http://nlp.stanford.edu/projects/dissertations/>

Chuang et al
CHI '12

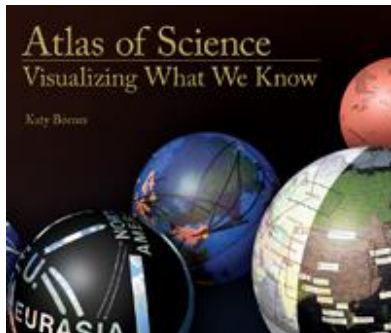
<http://nlp.stanford.edu/projects/dissertations/browser.html>

Fall 2016

CS 7450

61

Wonderful Book and Website



K. Börner



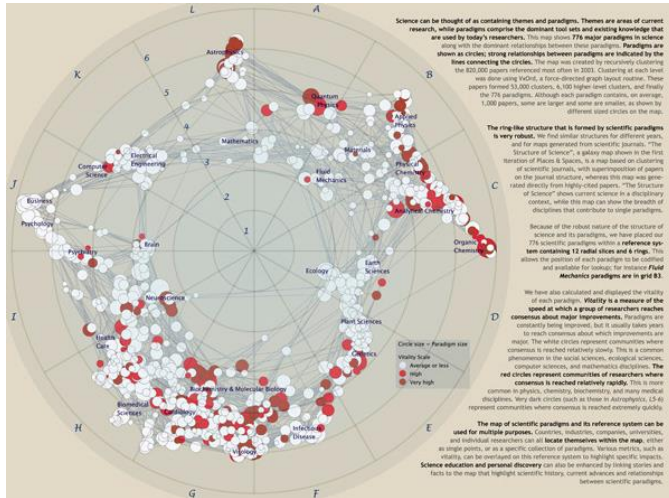
<http://scimaps.org>

Fall 2016

CS 7450

62

Some Examples



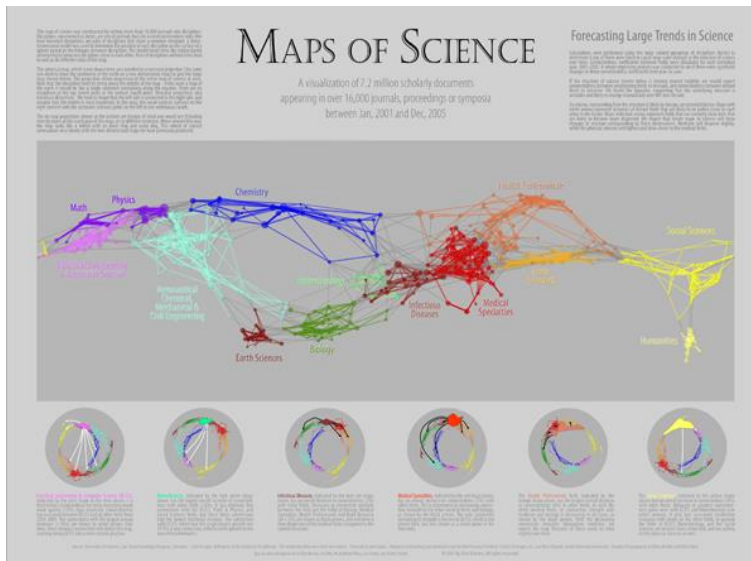
Boyack & Klavans

http://scimaps.org/maps/map/map_of_scientific_pa_55/

Fall 2016

CS 7450

63



Klavans & Boyack

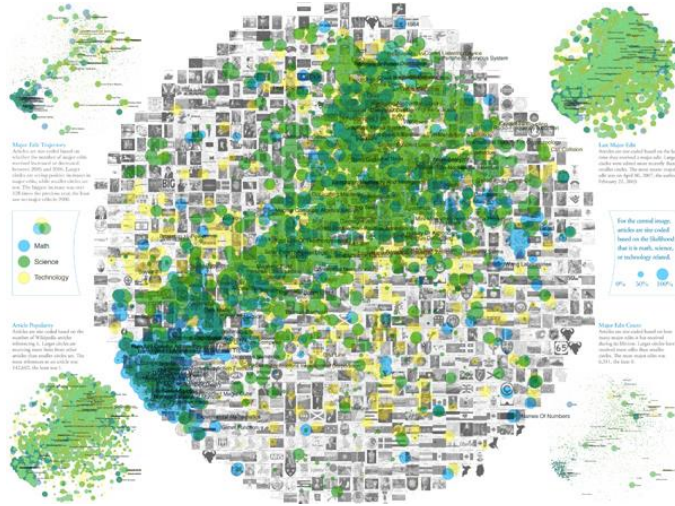
http://scimaps.org/maps/map/maps_of_science_fore_50/

Fall 2016

CS 7450

64

Science Related Wikipedia Activity



Allgood,
Herr,
Holloway &
Boyack

http://scimaps.org/maps/map/science_related_wiki_49/

Fall 2016

CS 7450

65

Temporal Issues

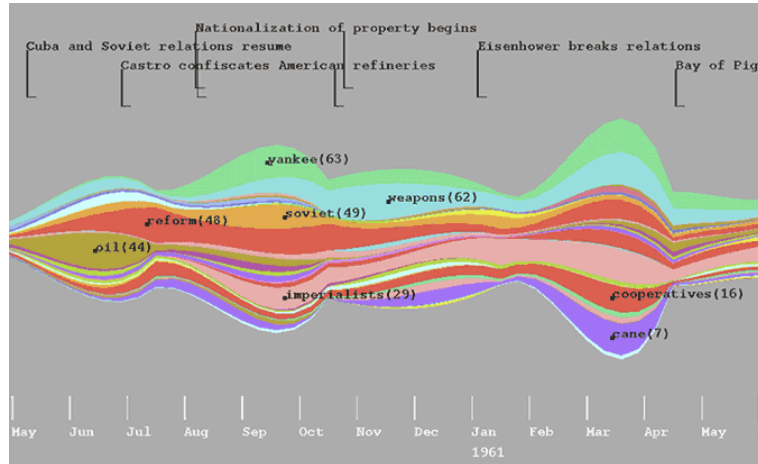
- Semantic map gives no indication of the chronology of documents
- Can we show themes and how they rise or fall over time?

Fall 2016

CS 7450

66

ThemeRiver



Havre, Hetzler, & Nowell
InfoVis '00

Fall 2016

CS 7450

67

Representation



- Time flows from left->right
- Each band/current is a topic or theme
- Width of band is "strength" of that topic in documents at that time

Fall 2016

CS 7450

68

More Information



- What's in the bands?
- Analysts may want to know about what each band is about

Topic Modeling



- Hot topic in text analysis and visualization
- Latent Dirichlet Allocation
- Unsupervised learning
- Produces "topics" evident throughout doc collection, each modeled by sets of words/terms
- Describes how each document contributes to each topic

TIARA



- Keeps basic ThemeRiver metaphor
- Embed word clouds into bands to tell more about what is in each
- Magnifier lens for getting more details
- Uses Latent Dirichlet Allocation to do text analysis and summarization

Liu et al
CIKM '09, KDD '10, VAST '10

Fall 2016

CS 7450

71

Representation

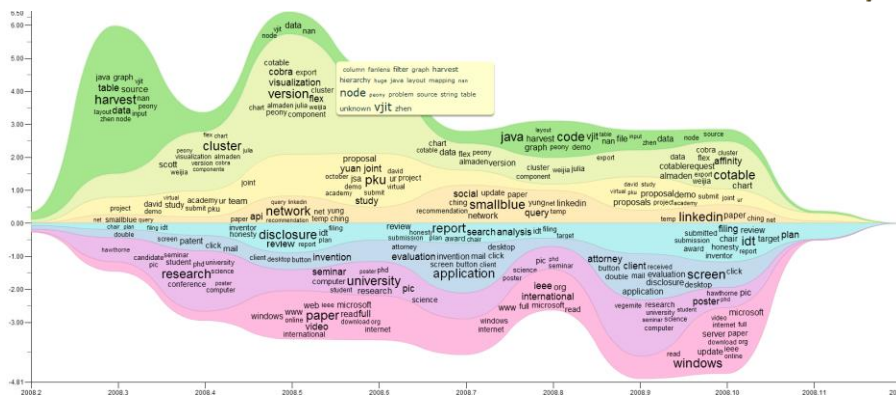


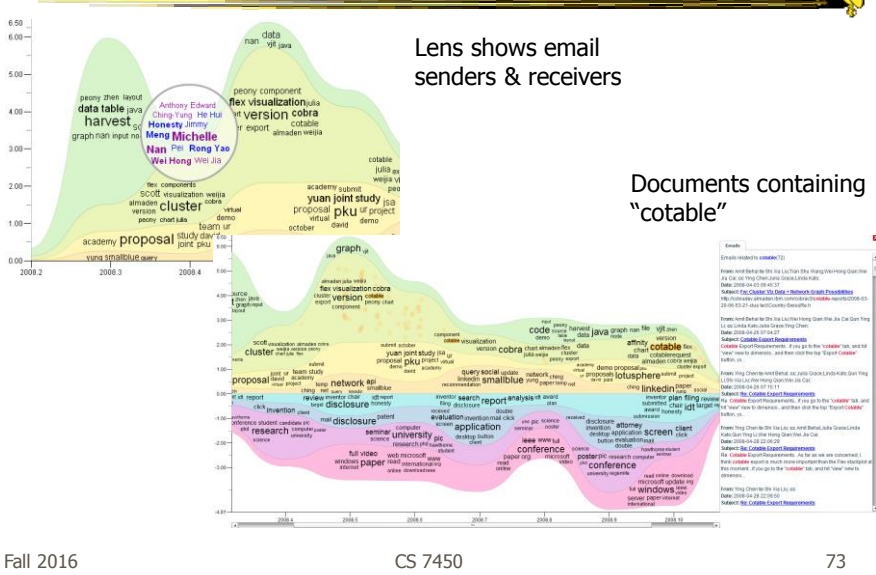
Figure 1. Annotated TIARA-created visual summary of 10,000 emails in the year 2008. Here, the x-axis encodes the time dimension, the y-axis encodes the importance of each topic. Each layer represents a topic, which is described by a set of keywords. These topic keywords are distributed along the time, summarizing the topic content and the content evolution over time. The tool tip shows the aggregated content of the top-most topic (green one).

Fall 2016

CS 7450

72

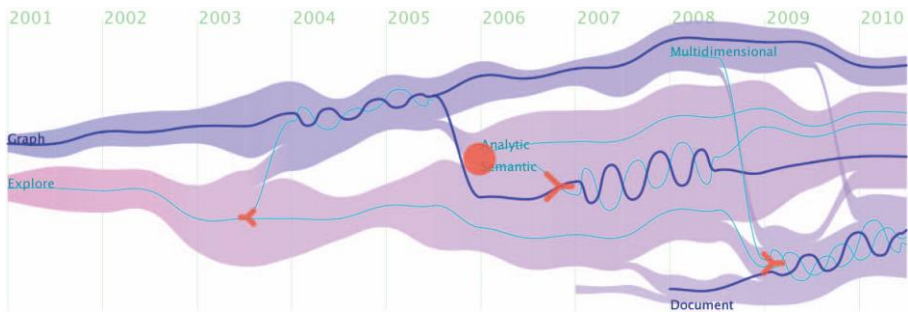
Features



Lens shows email senders & receivers

Documents containing "cotable"

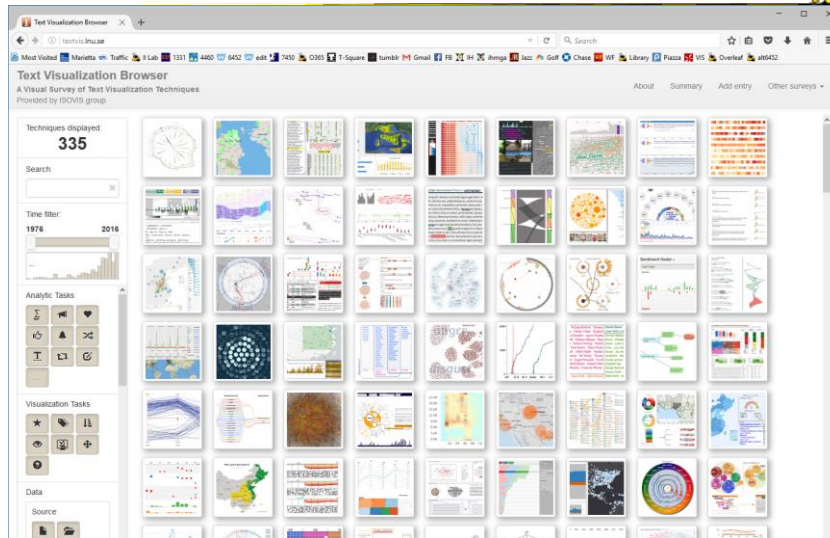
TextFlow



Showing how topics merge and split

Cui et al
TVCG (InfoVis) '11

Text Vis Browser



Fall 2016

CS 7450

75

Learning Objectives



- List different queries/tasks often needed on document collections
- List various analytic metrics often calculated on documents
- List different aspects of documents often visualized
- Explain vector space document analysis (similarity calculation, search)
- Explain TFIDF
- Describe visual representation used by and contributions of these systems
 - Themail, PaperLens, Jigsaw, ThemeScope/IN-SPIRE, ThemeRiver

Fall 2016

CS 7450

76

Reading



- Meirelles, chapter 6

Fall 2016

CS 7450

77

Upcoming



- Hierarchy and Tree data 1 & 2

Fall 2016

CS 7450

78