

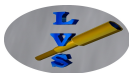


Introduction
ITMS
Preprocessing
Data
Data
Visualization
Cluster
Analysis
Topic
Modeling
Google Book
API
Future
Directions
References

Interactive Visual Data Analysis Part Two Interactive Text Mining Suite

Olga Scrivner

Indiana University



Workshop in Methods





Data Mining

Introduction

ITMS

Preprocessing
Data

Data
Visualization

Cluster
Analysis

Topic
Modeling

Google Book
API

Future
Directions

References

“As our collective knowledge continues to be digitized and stored (...) it becomes more difficult to find and discover what we are looking for.” (Blei 2012)



New Ways of Exploring Data Collections

Introduction

ITMS

Preprocessing
Data

Data
Visualization

Cluster
Analysis

Topic
Modeling

Google Book
API

Future
Directions

References

- Word clouds (Vuillemot et al., 2009)

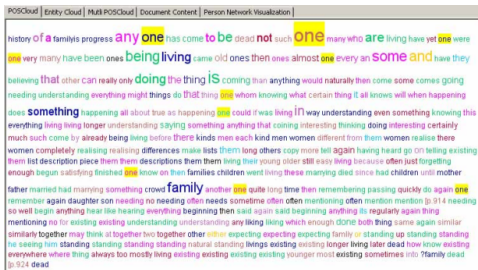


FIG. 4: Chapter 9, the word "one" highlighted in PosCloud visualization



Visualization Methods

Introduction

ITMS

Preprocessing
Data

Data
Visualization

Cluster
Analysis

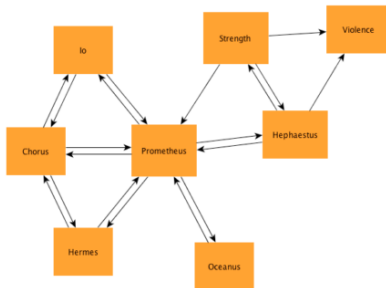
Topic
Modeling

Google Book
API

Future
Directions

References

- Social network graphs (Rydberg-Cox, 2011)





Visualization Methods

Introduction

ITMS

Preprocessing
Data

Data
Visualization

Cluster
Analysis

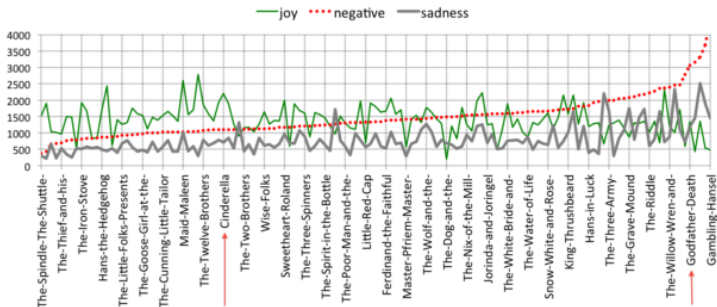
Topic
Modeling

Google Book
API

Future
Directions

References

- Tracking emotion and sentiment in fairy tales (Mohammad, 2012)





Topic Modeling

Introduction

ITMS

Preprocessing
Data

Data
Visualization

Cluster
Analysis

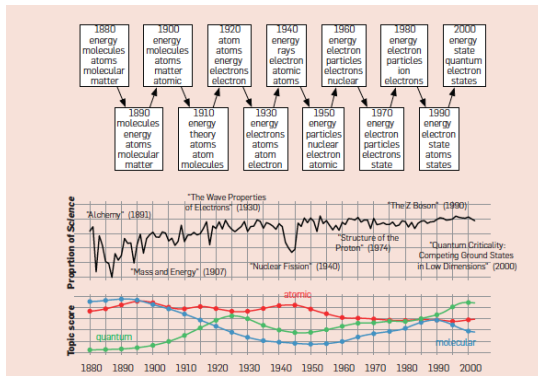
Topic
Modeling

Google Book
API

Future
Directions

References

Discovering underlying theme of collection from *Science* magazine
1990-2000 (Blei 2012)





Technological and Methodological Obstacles

Introduction

ITMS

Preprocessing Data

Data Visualization

Cluster Analysis

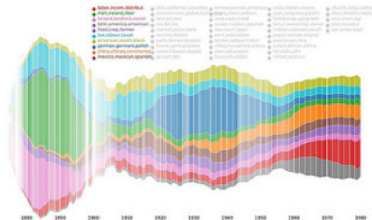
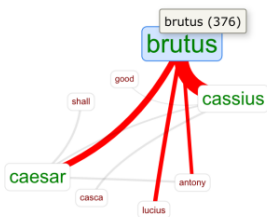
Topic Modeling

Google Book API

Future Directions

References

- Many tools require some programming skills (Mallet, Meta, R and Python libraries)
- GUI tools are limited to certain formats and functions (Voyant, PaperMachine)
- Lack of active control by users





Interactive Text Mining Suite

Introduction

ITMS

Preprocessing
Data

Data
Visualization

Cluster
Analysis

Topic
Modeling

Google Book
API

Future
Directions

References

- A user-friendly tool for quantitative analysis and visualization of unstructured data
- Platform-independent
- Interactive





ITMS Structure

Introduction

ITMS

Preprocessing
Data

Data
Visualization

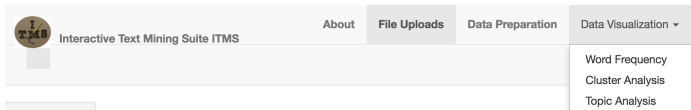
Cluster
Analysis

Topic
Modeling

Google Book
API

Future
Directions

References



1 File Uploads

- Upload files (txt, pdf, rdf and Google books API)

2 Data Preparation

- Data preprocessing (stopwords, stemming, metadata)

3 Data Visualization

- Word frequencies, Cluster analysis and topic modeling



ITMS Structure

Introduction

ITMS

Preprocessing
Data

Data
Visualization

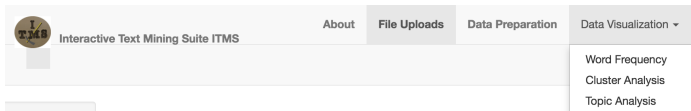
Cluster
Analysis

Topic
Modeling

Google Book
API

Future
Directions

References



1 File Uploads

- Upload files (txt, pdf, rdf and Google books API)

2 Data Preparation

- Data preprocessing (stopwords, stemming, metadata)

3 Data Visualization

- Word frequencies, Cluster analysis and topic modeling



Workshop Files

Introduction

ITMS

Preprocessing
Data

Data
Visualization

Cluster
Analysis

Topic
Modeling

Google Book
API

Future
Directions

References

- Download 3 text files

<http://ssrc.indiana.edu/seminars/wim.shtml>

- NY Times articles (3 documents in a plain text format)

- ITMS Web site:

<http://www.interactivetextminingsuite.com>





Upload File

Introduction

ITMS

Preprocessing
Data

Data
Visualization

Cluster
Analysis

Topic
Modeling

Google Book
API

Future
Directions

References

About

File Uploads

Data Preparation

Data Visualization ▾

Choose File(s) in TEXT format

Browse...

No file selected

Text Files

PDF Files

ZOTERO

Structured Data

POS-Tagged
Text



Upload File

Introduction

ITMS

Preprocessing
Data

Data
Visualization

Cluster
Analysis

Topic
Modeling

Google Book
API

Future
Directions

References

About

File Uploads

Data Preparation

Data Visualization ▾

Choose File(s) in TEXT format

Browse...

No file selected

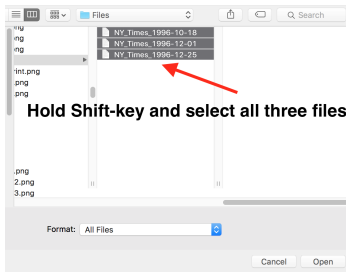
Text Files

PDF Files

ZOTERO

Structured Data

POS-Tagged
Text





Upload File

Introduction

ITMS

Preprocessing
Data

Data
Visualization

Cluster
Analysis

Topic
Modeling

Google Book
API

Future
Directions

References

About

File Uploads

Data Preparation

Data Visualization ▾

Choose File(s) in TEXT format

Browse...

No file selected

Text Files

PDF Files

ZOTERO

Structured Data

POS-Tagged
Text

Choose File(s) in TEXT format

Browse...

3 files

Upload complete

[1] "NY_Times_1996-10-18.txt"

"NY_Times_1996-12-01.txt" [3]

"NY_Times_1996-12-25.txt"

Corpus Size Total: 1611



Before performing data analysis we should preprocess data.

Metadata

The New York Times October 18, 1996, Friday, Late Edition Final Two Different Pleas for Change: Excerpts From a Second Senate Debate SECTION: Section B; Page 22; Column 1; Metropolitan Desk LENGTH: 1145 words Following are excerpts from last night's debate in Trenton between Robert G. Torricelli, a Democrat, and Richard A. Zimmer, a Republican, as transcribed by The New York Times. Opening Statements TORRICELLI For several months Dick Zimmer and I have been campaigning around New Jersey, asking for your help to get elected to the United States Senate. I know you're disappointed in the campaign. So am I. It's deteriorated into personal accusations and acrimony. But in truth, this campaign isn't about me. And it's not about Dick Zimmer. It's about you, your families and your future. It's also not about taxes or spending. I voted for a tax cut last year. So did Dick Zimmer. I voted for the



Preprocessing Options

Introduction

ITMS

Preprocessing
Data

Data
Visualization

Cluster
Analysis

Topic
Modeling

Google Book
API

Future
Directions

References

Select preprocessing options and click **apply**.

Select Preprocessing Steps

☒ Remove Punctuation

Exceptions (keep hyphen or apostrophe)

☒ none

☐ both

☐ hyphen

☒ Lower Case

☒ Remove Numbers

Preprocessing Viewer

Apply Steps or Default (no preprocessing)

☒ apply

☐ default



Stopwords

Introduction
ITMS

Preprocessing
Data

Data
Visualization

Cluster
Analysis

Topic
Modeling

Google Book
API

Future
Directions

References

Stopwords (e.g. **the**, **and**): select **Default** for English

Data Cleaning

Stopwords

Stemming

Metadata

Select Default or Upload

- ☐ None
- ☒ Default
- ☐ Upload

Default is the list from tm package:
stopwords("SMART")

[1] "a"	"a's"
[5] "above"	"accordir
[9] "actually"	"after"
[13] "against"	"ain't"
[17] "allows"	"almost"
[21] "already"	"also"
[25] "am"	"among"
[29] "and"	"another"
[33] "anyhow"	"anyone"
[37] "anyways"	"anywhere
[41] "appreciate"	"appropri
[45] "around"	"as"
[49] "asking"	"associat
[53] "away"	"awfully"



Manual Removal of Stopwords

Introduction

ITMS

Preprocessing
Data

Data
Visualization

Cluster
Analysis

Topic
Modeling

Google Book
API

Future
Directions

References

Based on the need, remove any additional stopwords that you may consider a noise, e.g, **paper**, **shows** etc

Manual Removal

Select one or multiple words (hold shift key down)

Select words to be removed

made written |
subject
supported
systems
textbooks
training
ultimately
union
voluntary
judiciary jurisprudence led legislative made
marriage member needed notes operation
organization parent perspective polity
possess practice preference procedures
produce progeny proper provided quality

Viewer

Apply Stopwords or None (no changes)

☒ apply

☐ none

private law field legislative intervention
rendering eu law part national legal system
requiring courts account jurisprudence
european association quality european organization

Select **apply**



Stemming

Introduction

ITMS

Preprocessing
Data

Data
Visualization

Cluster
Analysis

Topic
Modeling

Google Book
API

Future
Directions

References

To improve analytics, you can stem all your tokens, ex. instead of **worked**, **works**, **working**, you will have only one relevant stem **work**

Stems - tm package



Choose Language

- ☐ none
- ☒ English
- ☐ Spanish
- ☐ Danish
- ☐ Dutch
- ☐ Finnish

Stem Viewer

privat law field legis intervent render eu law part nation legal system requir court account jurisprud european societi politi proper organ oper legal system law appli studi sourc led reconfigur common law legal famili parent legal system enter marriag give rise progeni privat european communiti act european union law legal effect nation legal system basi nation court requir appli eu law remain conceptu strike heart domest legal system hold state court subject eu law requir note prefer australia canada zealand legal system analysi support earliar studi year train need qualifi practic favour legal system possess subconsci bias system procedur vis vis member state embed legal system conceiv state complianc eu law voluntari act reli extent echr case law part uk legal system dealt textbook uk academ general access english ultim produc judgment higher qualiti give judiciari perspect legal system court benefit fulli insight provid compar law



Metadata Extraction

Introduction

ITMS

Preprocessing
Data

Data
Visualization

Cluster
Analysis

Topic
Modeling

Google Book
API

Future
Directions

References

You can extract or upload metadata. You will need datestamp (year) information for chronological topic modeling.



Choose metadata source

- ☐ None
- ☐ From metadata of each uploaded PDF
- ☐ From separate CSV file
- ☐ From separate JSON file
- ☐ From separate XML file
- ☒ From zotero files metadata

Show 25 entries

Search:

date	title	author
2015	Gilliker - 2015 - The Influence of Eu and European Human Rights Law .pdf	Gilliker
<input type="text"/>	<input type="text"/>	<input type="text"/>

Showing 1 to 1 of 1 entries

Previous 1 Next



Visualization

Introduction

ITMS

Preprocessing
Data

Data
Visualization

Cluster
Analysis

Topic
Modeling

Google Book
API

Future
Directions

References

File Uploads

Data Preparation

Data Visualization

Frequency Table

Word Clouds

Length

KWIC

Punctuation

Data Visualization ▾

Word Frequency

Cluster Analysis

Topic Analysis



◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡



☐ Sans Serif

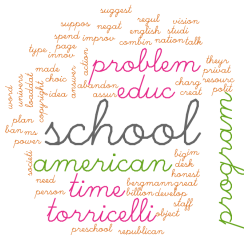
☒ Script

☐ Gothic

☐ black

☐ green

☒ multi





Cluster Analysis

File Uploads

Data Preparation

Data Visualization

Word Frequency

Cluster Analysis

Topic Analysis

You need to have at least **three** documents

Documents will be grouped based on their term similarity measures

Agglomeration Methods

Select method for cluster groups

- ☒ ward.D
- ☐ single
- ☐ complete
- ☐ average
- ☐ median
- ☐ centroid

Distance Measure

Select measure type

- ☒ euclidean
- ☐ maximum
- ☐ manhattan
- ☐ minkowski
- ☐ canberra
- ☐ binary



Cluster Analysis

Introduction

ITMS

Preprocessing
Data

Data
Visualization

Cluster
Analysis

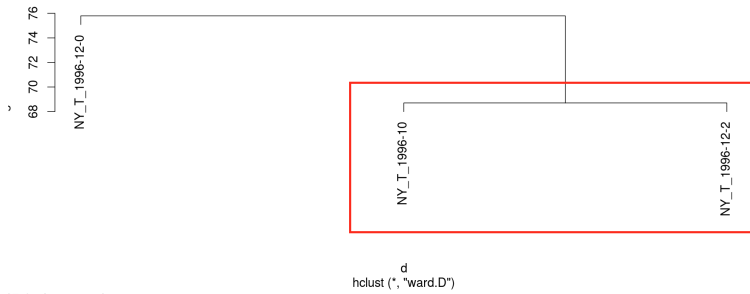
Topic
Modeling

Google Book
API

Future
Directions

References

Cluster Dendrogram





Topic Modeling

Introduction

ITMS

Preprocessing
Data

Data
Visualization

Cluster
Analysis

Topic
Modeling

Google Book
API

Future
Directions

References

- **LDA** (Latent Dirichlet allocation)
- **STM** (Structural Topic model)
- Chronological topic visualization (lda): requires metadata



Topic Modeling Tuning

- Introduction
- ITMS
- Preprocessing Data
- Data Visualization
- Cluster Analysis
- Topic Modeling
- Google Book API
- Future Directions
- References

- Selection of topics (how many different themes)
- Selection of words per theme (how many words per topic)
- Selection of iteration



Topic Model Selection

Introduction

ITMS

Preprocessing
Data

Data
Visualization

Cluster
Analysis

Topic
Modeling

Google Book
API

Future
Directions

References

File Uploads

Data Preparation

Data Visualization

Word Frequency

Cluster Analysis

Topic Analysis

Model Creation

LDA Visualization

STM Visualization

Metadata Topic
Visualization

Topic selection

Select number of topics - an integer representing the number of in the model. Default is 3.

Select or Type Number of Topics

3

Select the top number of words associated with a given topic. Default is 3.

Select or Type Number of Words per Topic

3



LDA Topic Model

Introduction

ITMS

Preprocessing
Data

Data
Visualization

Cluster
Analysis

Topic
Modeling

Google Book
API

Future
Directions

References

Model Creation

LDA Visualization

STM Visualization

Metadata Topic
Visualization

Run LDA Analysis

☐ none

☒ run

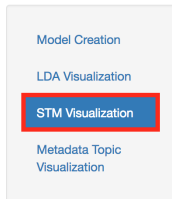
Selected Topics LDA (lda.collapsed.gibbs
package)

V1	V2	V3
policy	children	public
evidence	care	zimmer
president	vouchers	schools



STM Topic Model

- Introduction
- ITMS
- Preprocessing Data
- Data Visualization
- Cluster Analysis
- Topic Modeling
- Google Book API
- Future Directions
- References



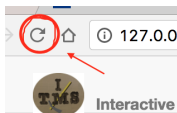
Structural Topics STM

Topic 1: compelling, definition, physical, choice, decisionmakers, influenced, notion, relationships, shifts, sweeping, trendthe, worknow, concepts, explained, made, behaviour, grows, practice, readings, collection
Topic 2: changing, covers, early, flexibility, key, magic, organizational, paradigmatic, technology, theories, applies, conceptualize, consumer, internal, investigation, models, perspective, rainbow, topics, vital
Topic 3: investigating, adaptation, analyze, barriers, institution, positions, yale, hundred, lens, procedure, capture, central, concern, fixed, notions, shaping, sketches, audience, balance, dispute

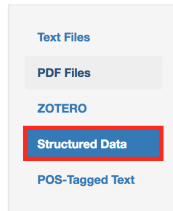
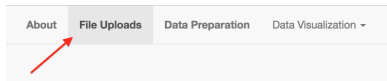


Other Formats - Google Books

Before switching to other data formats, refresh your local browser.



Start with **File Uploads** and select **Structured Data**





Other Formats - Google Books

Introduction

ITMS

Preprocessing
Data

Data
Visualization

Cluster
Analysis

Topic
Modeling

Google Book
API

Future
Directions

References

Select your search terms and submit

Choose file format

☐ XML

☐ JSON

☒ Google Books Search

Enter your search terms for Google Books,
separated by spaces

social science

Submit

Current limitation is 40 books

Show 25 entries

Search:

titles	authors	dates	corpus
Readings in the Philosophy of Social Science	Michael Martin, Lee C. McIntyre	1994	Readings in the Philosophy of Social Science the first comprehensive anthology in the



Future Options

Shiny Web Application is highly customizable

- 1 Part-of-speech tagging (tm package)
- 2 Network analysis (igraph package)
- 3 Name Entity Recognition (NLP package)
- 4 Twitter Streaming (twitterR package) - will requires user's twitter set-up for streaming but information will be provided how to set it up

Open for other suggestions and collaboration - contact
obscrivn@indiana.edu



Acknowledgements

Introduction

ITMS

Preprocessing
Data

Data
Visualization

Cluster
Analysis

Topic
Modeling

Google Book
API

Future
Directions

References

I would like to thank WIM for providing this opportunity.

Contributors: Jefferson Davis, Irina Trapido, Jay Lee



References I

- [1] Many open source R packages: tm, shiny, NLP, stringi, stringr, topicmodels, lda and many more
- [2] Baayen, Harald. 2008. *Analyzing linguistic data: A practical introduction to statistics*. Cambridge: Cambridge University Press
- [3] Gries, Stefan Th. 2015. *Quantitative designs and statistical techniques*. In Douglas Biber Randi Reppen (eds.), *The Cambridge Handbook of English Corpus Linguistics*. Cambridge: Cambridge University Press
- [4] Jockers, Matthew. 2014. *Text Analysis with R for Students of Literature. Quantitative Methods in the Humanities and Social Sciences*. Springer International Publishing, Cham
- [5] Moretti, Franco. 2005. *Graphs, Maps, Trees: Abstract Models for a Literary History*. Verso
- [6] Oelke, Daniella, Dimitrios Kokkinakis, and Mats Malm. 2012. Advanced visual analytics methods for literature analysis. *Proceedings of the 6th EACL Workshop on Language Technology for Cultural Heritage, Social 561Sciences, and Humanities*, pages 35–44
image credits: <https://media.giphy.com/media/10zsjaH4g0GgmY/giphy.gif>