

Computational Text Analysis for Content Analysis

Taught By Sara Morrell and Mel Williams
Digital Integration Teaching Initiative (DITI)

POLS 7346 Resilient Cities
Daniel Aldrich
Spring 2026

Workshop Agenda

- Introduction to key terms and concepts in computational text analysis (CTA).
- Discussion of CTA's applications and uses in research.
- Introduction to web-based text analysis tools.
 - Word Counter, Word Trees, Voyant, Lexos

Slides: <https://bit.ly/sp26-aldrich-pols7346-text-analysis>

Data: <https://bit.ly/sp26-aldrich-pols7346-text-analysis-data>

What is Computational Text Analysis?

Feel free to ask questions at any point during the presentation!

Computational Text Analysis

Computational text analysis refers to the **array of methods used to “read” texts with a computer**. It is similar to statistical analysis, but the data is texts (words) instead of numbers.

Text analysis:

- Involves a computer drawing out patterns in a text, and a researcher interpreting those patterns.
- Includes methods such as word count frequency, keywords in context, computational modeling (with machine learning), and sentiment analysis.
- Is conducted using web-based tools or coding languages like Python and R.

Why Computational Text Analysis?

Computational text analysis can help us **analyze very large amounts of data**, **identify keywords**, and **discover patterns** in texts. Using text analysis, researchers may find surprising results that they would not have discovered from traditional methods alone.

From collections of texts, researchers can discover keywords that serve as a proxy for major trends in societies, cultures, and policies. For example, computational tools can reveal patterns on how public officials communicate policies, which issues are of concern, which phrases leaders regularly employ, and much more.

Language Used in Disaster News Coverage



Go to the [Television Explorer](#). Search “Hurricane” or other disaster-related terms.

- What do you notice about the TV coverage of these terms? What is surprising?
 - How might this language shape policies?

Feel free to ask questions at any point during the presentation!

Key Terms (1/2)

- **Corpus (plural-corpora):** A collection of texts used for analysis and research purposes.
- **Stop words:** Words that appear frequently in a language, like pronouns, prepositions, and basic verbs. These are often removed for computational analysis. Some English stop words include: a, the, she, he, I, me, us, of, is, would, could, should, etc.
- **Word Count Frequency:** Counting the total times a word appears in a text/corpus or the percentage of how often it appears.

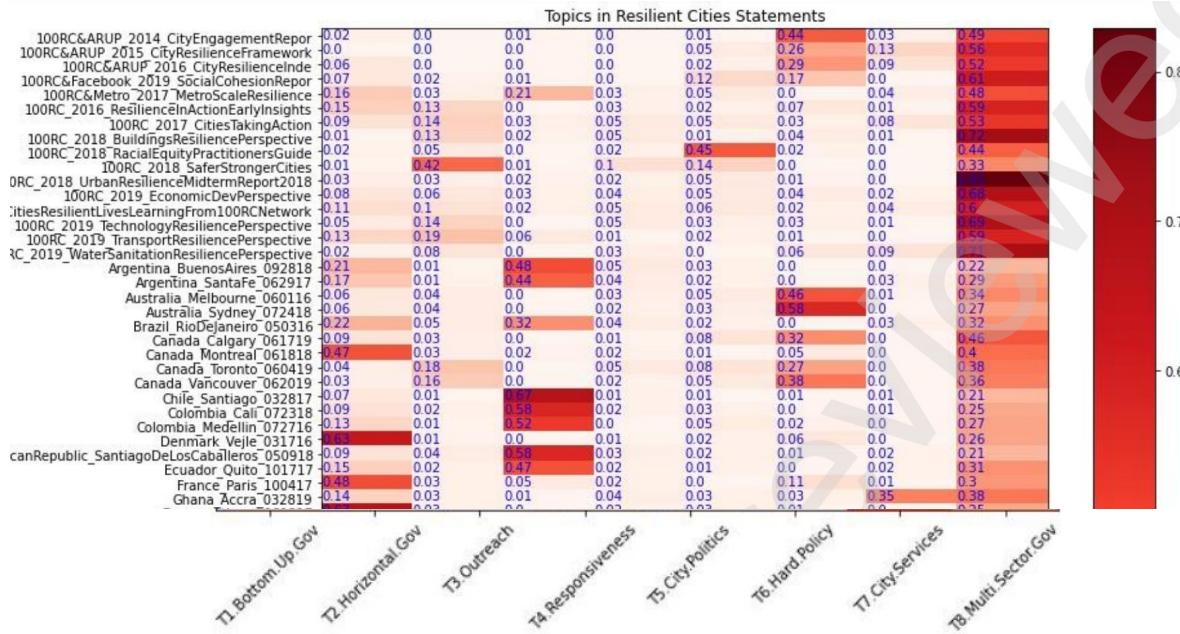
Key Terms (2/2)

- **nGram:** A continuous sequence of n items in a text. A bigram (or 2 continuous words) could be ‘United States,’ while a trigram (3 words) could be ‘yes we can.’
- **Sentiment Analysis:** Measuring the sentiment of a text based on a scale such as negative/positive or happy/sad. Each word has a particular weight to determine where on the scale it falls, and these weights are calculated to determine a text’s overall sentiment.

Examples from Practice

*Feel free to ask questions at any point
during the presentation!*

Resilient Cities Statements Computational Text Analysis



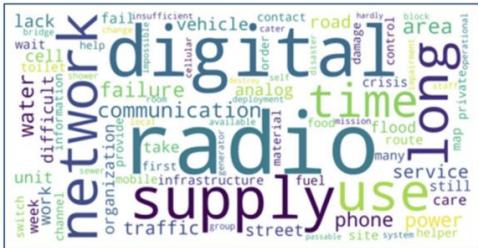
Computational text analysis of resilience strategy language shows different place-based priorities, as well as gaps and overlaps in 100 cities' resilience strategies.

DITI and POLS 7346 class alum Garrett Morrow applied computational text analysis to model and identify topics in resilient cities' strategy documents.

Disaster Resilience



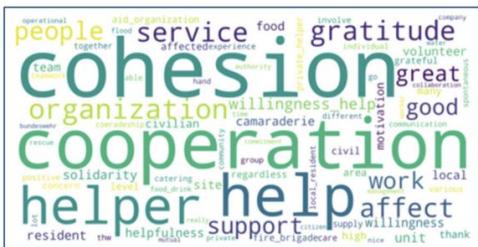
Survey question 1



Survey question 2



Survey question 3



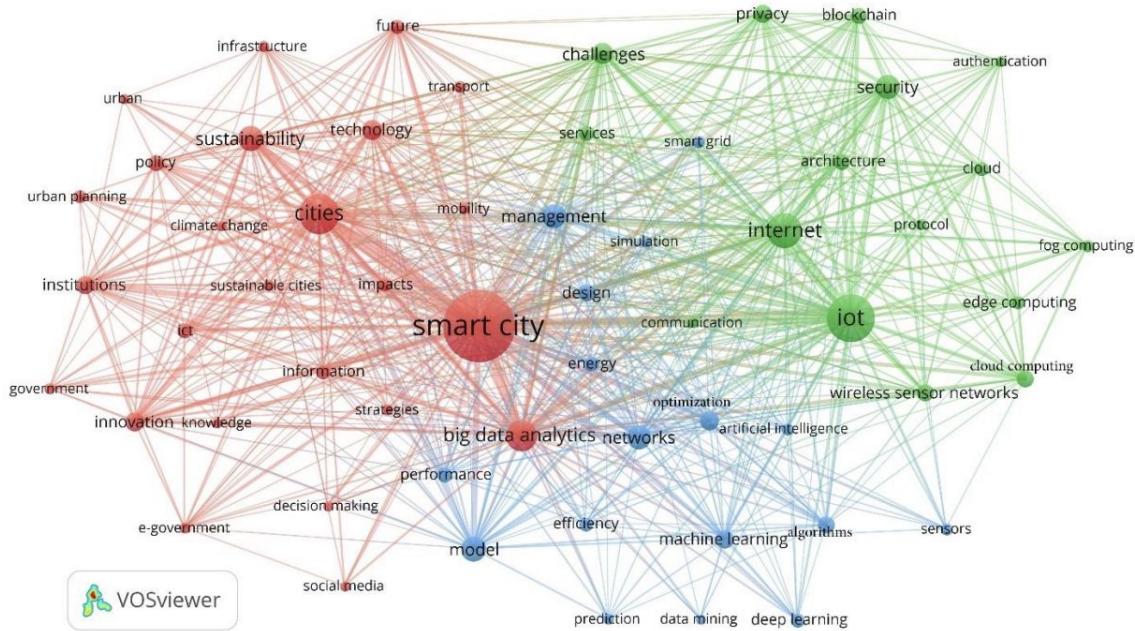
Survey question 4

“Word cloud of survey questions” (p. 4235)

Moghadas, M., Fekete, A., Rajabifard, A., & Kötter, T. (2023). The wisdom of crowds for improved disaster resilience: a near-real-time analysis of crowdsourced social media data on the 2021 flood in Germany. *GeoJournal*, 88(4), 4215-4241.

<https://doi.org/10.1007/s10708-023-10858-x>

Smart Cities



“The output of the term co-occurrence analysis for the second period (2016–2021)” (p. 11)

Sharifi, A., Allam, Z., Feizizadeh, B., & Ghamari, H. (2021). Three decades of research on smart cities: Mapping knowledge structure and trends. *Sustainability*, 13, 7140.

<https://doi.org/10.3390/su13137140>

Text Preparation

*Feel free to ask questions at any point
during the presentation!*

Corpus Building

Questions to consider as you begin your research:

- What are my research questions and why am I creating a corpus?
- What am I asking my corpus to do?
- What text(s) should form my corpus to answer my research questions?
- How should I organize my corpus to streamline my research processes and save time?

For more information, see our [Corpus Building Handout](#).

Preparing Your Text

1. Choose the texts or text selections that you would like to include.
2. Create a folder on your computer or cloud storage where you will store your corpus. Give it a clearly descriptive name, without spaces or special characters.
3. Copy and paste the text into a **plain text editor** (on Macs: Text Edit; on Windows: Notepad)
 - a. Mac users, you may need to make your Text Edit into a ‘plain text’. Open Text Edit, go to Preferences, and make sure “plain text” is selected
4. Save the text as a plain text file (with a .txt extension). Always make sure to name your files so you know what is in them!
5. Repeat steps above for each text in the corpus.

Our Text

We will use six political statements made in January 2025 in response to the LA urban fires:

- [Proclamation Of A State Of Emergency](#) by Gavin Newsom, Jan 7th, 2025
- [Statement from President Joe Biden](#), Jan 7th, 2025
- [Executive Order N-2-25](#) by Gavin Newsom, Jan 8th, 2025
- [Statement from President Joe Biden](#), Jan 13th, 2025
- [Inaugural Address](#) by Donald Trump, Jan 20th 2025
- [Putting People Over Fish: Stopping Radical Environmentalism to Provide Water to Southern California](#) by Donald Trump, Jan 20th, 2025

Sample Corpus

The sample .txt files are available on:

<https://bit.ly/sp26-aldrich-pols7346-text-analysis-data>

- You can download the files individual or click the Download all in the upper right corner
- If you download all, a zipped folder will download containing the files
 - On Mac: Double click the folder to unzip it
 - On PC: Right click the folder and select Extract all

Initial Corpus Analysis

Open any one of the texts from the sample corpus:

What can you observe about the text? How long is it? What kinds of language does it use? What kinds of analysis might you do with a text like this?

Scan through a few more: do they seem largely similar? What do you think might be different?

Exploratory Tool: Word Counter

*Feel free to ask questions at any point
during the presentation!*

Word Counter

- <https://databasic.io/en/wordcounter/>
- A user-friendly **word counting tool**
- Allows you to count words, bigrams, and trigrams in plain text files and to download spreadsheets with your results
- The max file upload is 10MB
- The default is to lowercase all words and apply stop words, but you can change those settings
- For more information, please see:
<https://bit.ly/handout-data-basics-suite>

Word Counter Example

What seems
significant in the
most frequent
terms from the
Proclamation of a

State of Emergency?



This is a **word cloud**, used to get a sense of the **most used words in a document**. Words used more often are bigger than those used less often.

Feel free to ask questions at any point during the presentation!

“Tokenizing” text

Before words can be counted, they must be “tokenized” or divided into components that programs can treat as distinct segments. Different programs will have different standards for tokenization—this one uses both white spaces and punctuation marks (such as commas) to separate words into tokens. **What are some limitations of this approach?**

Data preparation

Go to the [upload/paste screen for WordCounter](#) and un-click the “ignore stop words” and “ignore case” options, then upload the Proclamation Of A State Of Emergency and count the words again.

What happened? Why do you think the default is to ignore stop words and remove differences between upper/lowercase words? Can you think of any limitations to this approach?

Bigrams and Trigrams

TOP WORDS ↓		BIGRAMS ↓		TRIGRAMS ↓	
Word	Frequency	bigram ↗	Frequency	trigram ↗	Frequency
emergency	22	of the	18	the palisades fire	9
state	21	the palisades	10	government code section	8
code	15	government code	10	in los angeles	6
california	13	of emergency	9	palisades fire and	6
government	13	palisades fire	9	fire and windstorm	6
palisades	10	code section	9	and windstorm conditions	6
fire	10	state of	7	los angeles and	5
conditions	10	and whereas	7	angeles and ventura	5
section	10	windstorm conditions	7	and ventura counties	5

In addition to single words, it is also useful to consider **bigrams** and **trigrams** which include additional context.

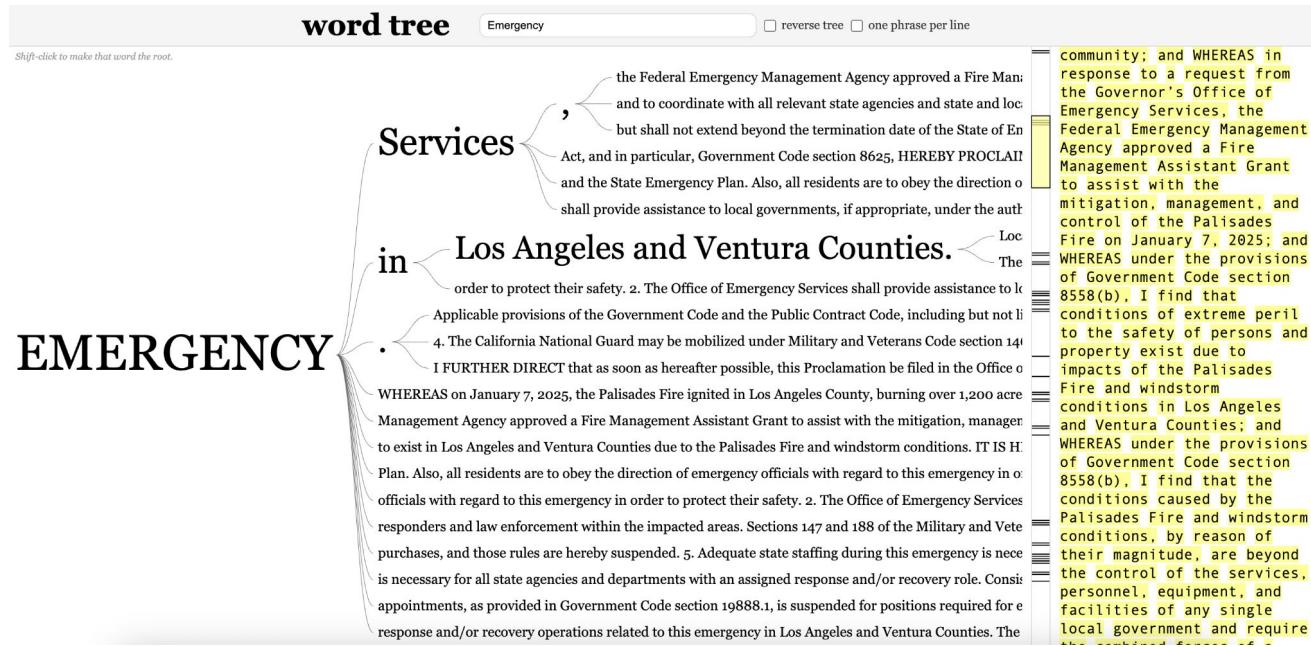
Exploratory Tool: Word Tree

*Feel free to ask questions at any point
during the presentation!*

Word Tree

- <https://www.jasondavies.com/wordtree/>
- A word tree **depicts multiple parallel sequences of words.**
- This is a good way to see patterns in word usage, based on words that appear before and after a term or terms of interest.
- There are some restrictions in size with this tool: fewer than 1 million words should work.
- Upload your text, enter a keyword or phrase to search, then try reversing the tree.
- It's often useful to search frequent terms identified by WordCounter

Word Tree Example



Word tree starting
with “Emergency”
from the
Proclamation Of A
State Of Emergency

What other terms
should we try?

*Feel free to ask questions at any point
during the presentation!*

Tools for corpus exploration: Voyant

*Feel free to ask questions at any point
during the presentation!*

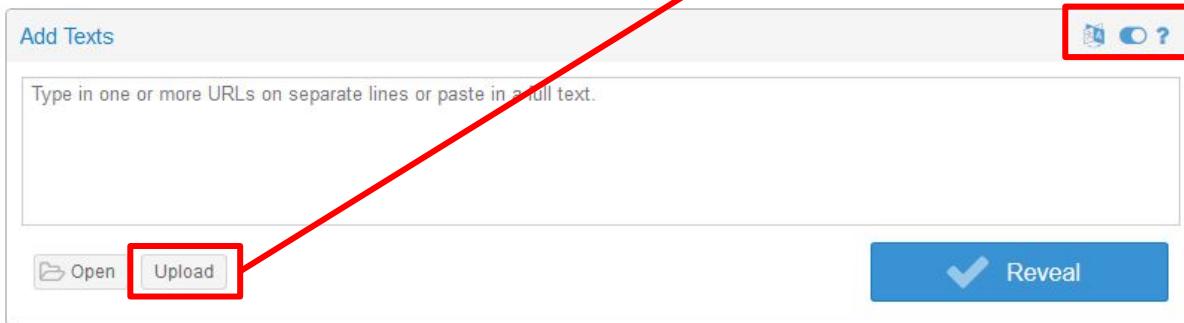
Voyant

Voyant makes it possible to **perform analyses on one or multiple files in many ways**, including word counts, nGrams (n=number of words), word frequency distributions, word trends across documents, and concordances.

<https://voyant-tools.org/>

For more information, see: **<https://bit.ly/handout-voyant-intro>**

Voyant: Upload



Click on Upload and navigate to the folder with the text documents you wish to analyze.

Alternatively, insert URLs or full text into the textbox.

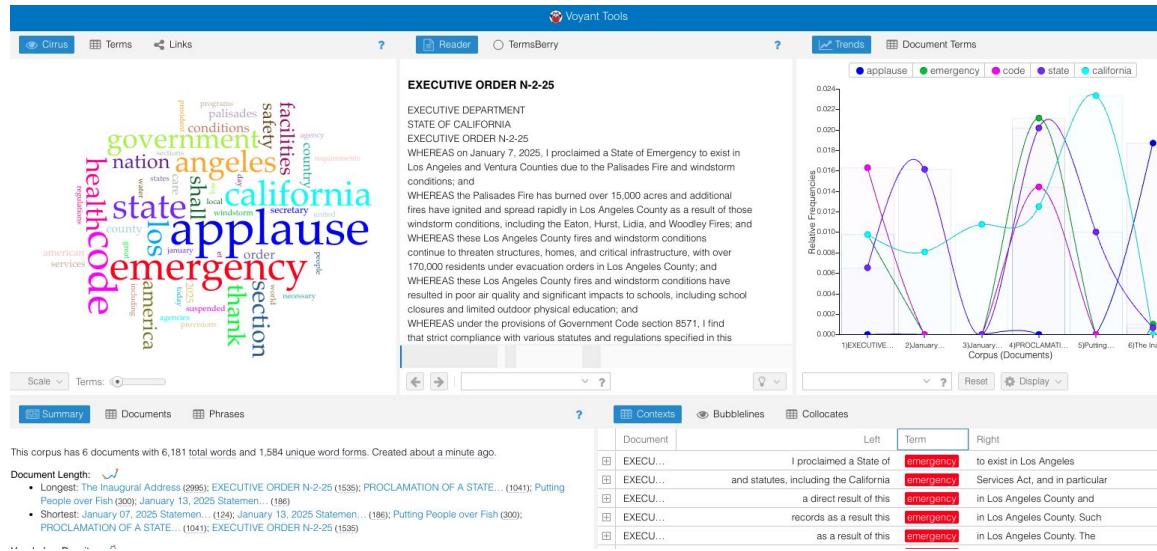
Click here for help and advanced options

Voyant: Dashboard

Results:

After you upload your corpus, you will see the default results page with multiple panes:

- A word cloud
- Reader section
- Trends
- Document Summary
- Word Contexts



These boxes can all be changed!

Feel free to ask questions at any point during the presentation!

Voyant: Changing Displayed Results

Hover on the right top corner of a pane and buttons will appear. Select the panes button and choose a new option from the dropdown menu. For example, we might want to try out the "Collocates" tool instead of the word cloud. Click on the '?' to learn more about how the tool works.

The image shows the Voyant interface. On the left, there is a word cloud visualization with words like 'applause', 'california', 'emergency', 'thank', 'state', 'los', 'applause', 'government', 'health', 'Code', 'nation', 'angels', 'shall', 'local', 'windstorm', 'secretary', 'united', 'people', 'try', 'order', 'et', 'water', 'care', 'states', 'county', 'american', 'services', 'great', 'january', 'today', '2025', 'suspended', 'provisions', 'agencies', 'including', 'necessary', 'regulations'. On the right, there is a sidebar menu with options: 'Terms' (selected), 'Links', 'Corpus Tools', 'Document Tools', 'Visualization Tools', 'Grid Tools', and 'Other Tools'. A red box highlights the 'Corpus Tools' button, and a red arrow points from it to a dropdown menu. The dropdown menu has 'Terms' (selected), 'Links', 'Collocates' (highlighted in blue), and a question mark icon. Below the sidebar is a table titled 'Collocates' with columns 'Term', 'Collocate', and 'Count (context)'. The table lists various collocates with their counts: applause (thank, 30), los (angeles, 30), code (section, 25), code (health, 15), los (county, 15), angeles (county, 15), health (safety, 15), health (code, 14), state (california, 13), california (state, 13), emergency (los, 12), emergency (angeles, 12), los (emergency, 12), angeles (emergency, 12), government (code, 12). At the bottom of the table are buttons for 'Scale' and 'Context' set to 935.

Term	Collocate	Count (context)
applause	thank	30
los	angeles	30
code	section	25
code	health	15
los	county	15
angeles	county	15
health	safety	15
health	code	14
state	california	13
california	state	13
emergency	los	12
emergency	angeles	12
los	emergency	12
angeles	emergency	12
government	code	12

Feel free to ask questions at any point during the presentation!

Voyant: Contexts (concordances)

Contexts, or concordances, show the different contexts around particular search terms. For example, you can see all the times the word “emergency” appears in the text and the contexts in which it appears.

The screenshot shows the Voyant interface with the 'Contexts' tab selected. A search bar at the bottom left contains the word 'emergency'. Below the search bar is a table with four columns: 'Document', 'Left', 'Term', and 'Right'. The 'Term' column is highlighted in red. The table lists several instances of the word 'emergency' found in the document, with context snippets from the text to its left and right. A red box highlights the search input field and the first few rows of the table.

Document	Left	Term	Right
EXECU...	I proclaimed a State of	emergency	to exist in Los Angeles
EXECU...	and statutes, including the California	emergency	Services Act, and in particular
EXECU...	a direct result of this	emergency	in Los Angeles County and
EXECU...	records as a result this	emergency	in Los Angeles County. Such
EXECU...	as a result of this	emergency	in Los Angeles County. The
EXECU...	as a result of this	emergency	in Los Angeles County. Such
EXECU...	cleanup of debris from this	emergency	or to address other impacts
EXECU...	of the effects of the	emergency	, or conducting other fire recovery

What other terms should we try?

Voyant: Tools for further exploration

- Voyant's [Getting Started](#) guide
- Voyant's [List of Tools](#), showing all the features possible with Voyant including descriptions of each
- Some useful tools to explore:
 - MicroSearch
 - Topics
 - Correlations
 - Collocates Graph

Tools for corpus exploration: Lexos

*Feel free to ask questions at any point
during the presentation!*

Lexos

Lexos provides a step-by-step guide for text uploading, preparation, and analysis.

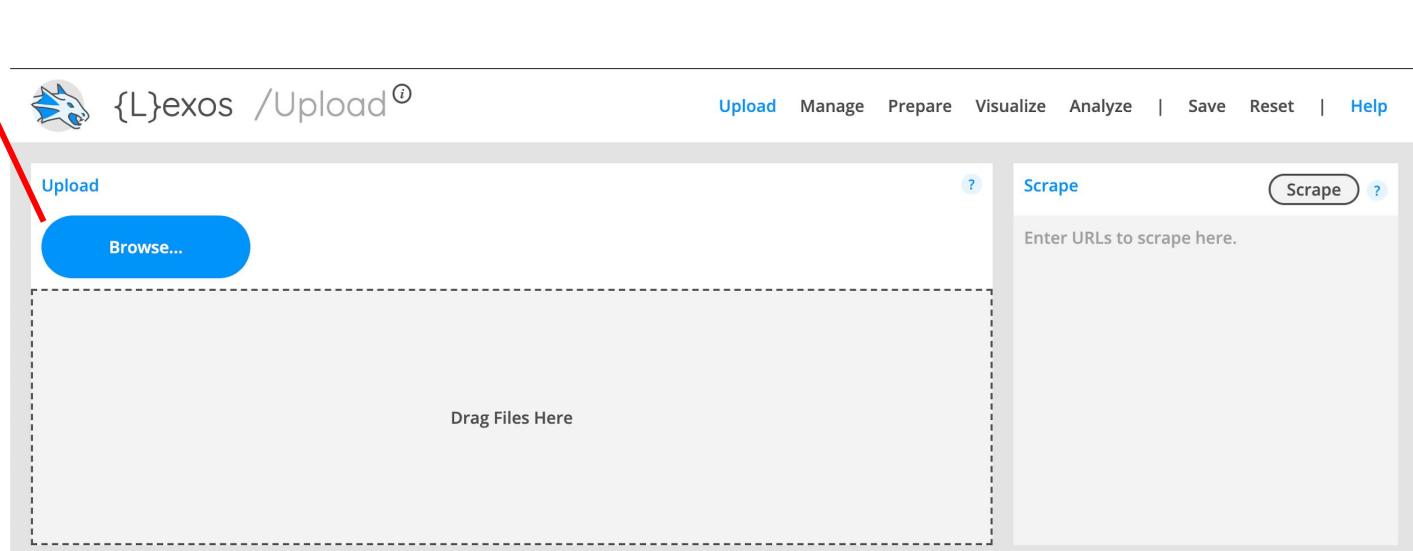
- **Upload:** upload your .txt file
- **Manage:** select the files you want to prepare and analyze
- **Prepare:** prepare your text for analysis
- **Visualize:** create visualizations of patterns across your corpus or in single texts
- **Analyze:** analyze your text

<http://lexos.wheatoncollege.edu/upload>

For more information, please see: <https://bit.ly/handout-Lexos-intro>

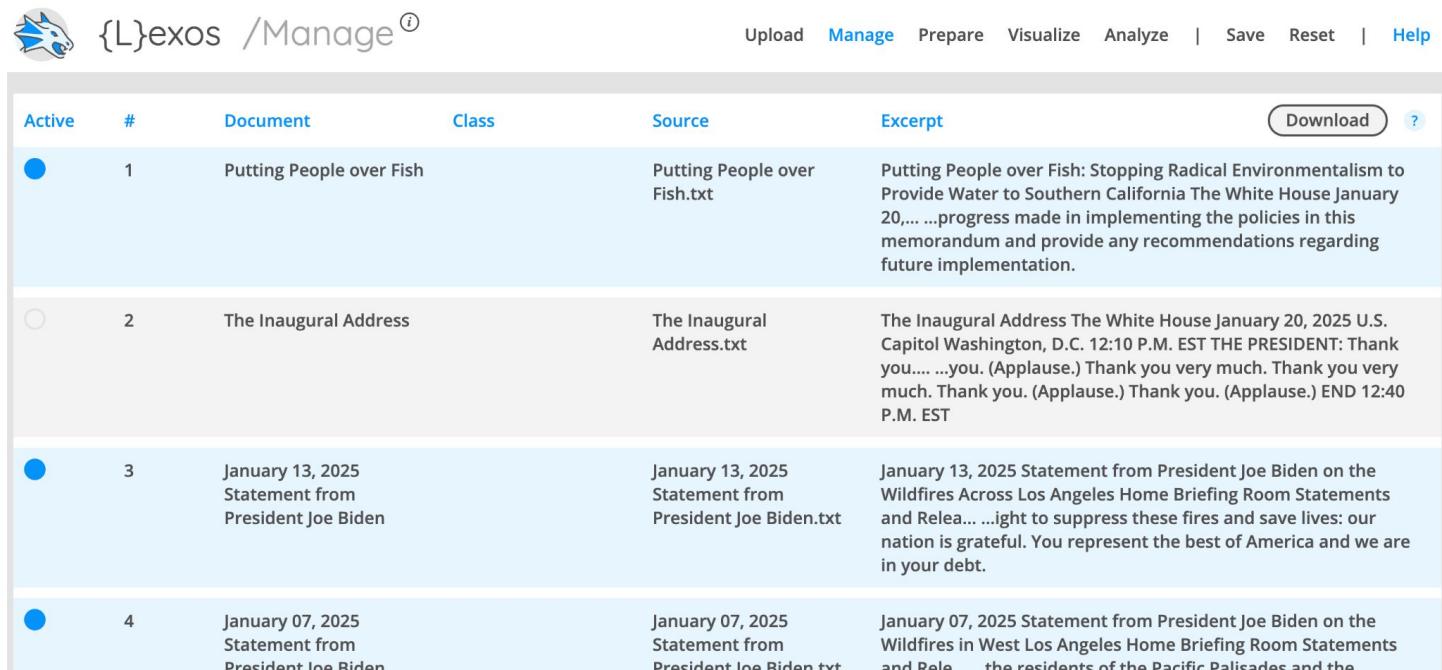
Lexos: Upload

Click Browse and select your entire text (or drag file into the “Drag Files Here” area). It can be easy to miss when the upload is done—click “Manage” to double check that the text file is there.



Lexos: Manage

Make sure the document you want to use is selected (blue = selected, gray = not selected)



The screenshot shows the Lexos Manage interface. At the top, there is a logo of a blue dragon-like creature, the word '{L}exos', and a '/Manage' button with an 'i' icon. To the right are navigation links: Upload, Manage (which is highlighted in blue), Prepare, Visualize, Analyze, and buttons for Save, Reset, and Help. Below this is a table with the following columns: Active, #, Document, Class, Source, and Excerpt. There are five rows of data:

Active	#	Document	Class	Source	Excerpt
<input checked="" type="radio"/>	1	Putting People over Fish		Putting People over Fish.txt	Putting People over Fish: Stopping Radical Environmentalism to Provide Water to Southern California The White House January 20,... ...progress made in implementing the policies in this memorandum and provide any recommendations regarding future implementation.
<input type="radio"/>	2	The Inaugural Address		The Inaugural Address.txt	The Inaugural Address The White House January 20, 2025 U.S. Capitol Washington, D.C. 12:10 P.M. EST THE PRESIDENT: Thank you.... ...you. (Applause.) Thank you very much. Thank you very much. Thank you. (Applause.) Thank you. (Applause.) END 12:40 P.M. EST
<input checked="" type="radio"/>	3	January 13, 2025 Statement from President Joe Biden		January 13, 2025 Statement from President Joe Biden.txt	January 13, 2025 Statement from President Joe Biden on the Wildfires Across Los Angeles Home Briefing Room Statements and Relea... ...ight to suppress these fires and save lives: our nation is grateful. You represent the best of America and we are in your debt.
<input checked="" type="radio"/>	4	January 07, 2025 Statement from President Joe Biden		January 07, 2025 Statement from President Joe Biden.txt	January 07, 2025 Statement from President Joe Biden on the Wildfires in West Los Angeles Home Briefing Room Statements and Relea... ...the residents of the Pacific Palisades and the

Lexos: Prepare (Scrub Case and Punctuation)

Lexos demonstrates some more advanced options you have for preparing your corpus. By “scrubbing,” you are transforming the texts in your corpus and making choices that will impact your results. Here are some possibilities:

- **Make Lowercase:** make all your letters lowercase. Even though you know “A” and “a” are the same letter, the computer treats these as two separate characters. Lowercasing removes this distinction.
- **Remove Punctuation:** remove punctuation, which may influence your results.

Lexos: Prepare (Scrub Words)

You can also stem words and remove certain words. Here are some possibilities:

- **Stop/Keep Words:** remove a list of words. Usually these would be **stop words**. With WordCounter, you had to use the stop words list the tool provided—now, you can choose your own.
- **Lemmas:** standardize to the *stem* of word. For example, you can stem all forms of the verb talk: talking, talked, talks, etc. to “talk”

Lexos: Removing Stop Words

Get a list of English stop words here:

<https://gist.github.com/sebleier/554280>

Copy and paste the stop words (hit "raw", then select all and copy) into the "Stop/Keep Words" box then select "Stop".

The screenshot shows the Lexos Scrub interface. In the top right, there are tabs: Upload, Manage, **Prepare**, Visualize, Analyze, Save, Reset, and Help. Below these are sections for Scrubbing Options, Stop/Keep Words, Previews, and Lemmas.

- Scrubbing Options:** Includes checkboxes for Make Lowercase, Remove Digits, Remove Spaces, Remove Tabs, Remove Newlines, Scrub Tags, Remove Punctuation, Keep Hyphens, Keep Apostrophes, and Keep Ampersands.
- Stop/Keep Words:** A section where users can choose to Off, Stop, or Keep stop words. The "Stop" option is selected. A red box highlights this selection. Below it is a list of stop words: applause, i, me, my, myself, we, our, ours, and a period.
- Previews:** Shows two examples:
 - Putting People over Fish:** putting people fish stopping radical environmentalism provide water southern california white house january january memorandum... ...r shall report regarding progress made implementing policies memorandum provide recommendations regarding future implementation
 - The Inaugural Address:** inaugural address white house january us capitol washington dc pm est president thank thank much everybody wow thank much... ...d way americans future golden age begun thank god bless america thank thank much thank much thank thank end pm est
- Lemmas:** A section for entering lemmas with an Upload button.
- Consolidations:** A section for entering consolidations with an Upload button.
- Special Characters:** Options for None, MUFI 3, MUFI 4, Early English HTML, and Old English SGML. The "None" option is selected.

You can also add stop words particular to your corpus.

Lexos: Applying your Preparations

BEFORE PREP

Previews

Preview

Apply

Download

Putting People over Fish

Putting People over Fish: Stopping Radical Environmentalism to Provide Water to Southern California
The White House January 20,... ...progress made in implementing the policies in this memorandum and provide any recommendations regarding future implementation.

AFTER PREP

Previews

Preview

Apply

Download

Putting People over Fish

putting people fish stopping radical environmentalism provide water southern california white house january january memorand...r shall report regarding progress made implementing policies memorandum provide recommendations regarding future implementation

Once you have made decisions about your preparations, click “**Apply**” and wait a few minutes. Because the program is going through each document and completing all the processes you selected, it needs some time. Then, you will see the final results of your preparation! You can also **download** your new corpus.



Lexos: Analyze > Top Words

The top words tool lets you compare word usage between individual documents and your corpus as a whole. If you want to make more specific comparisons, you can also assign “classes” to subsets of tools with the “Manage” screen.

- Words with high positive scores are **used more often** in each document, relative to the rest of the corpus.
- Words with high negative scores are **used less often**.

Hit the “Generate” button to see the top words for your texts.

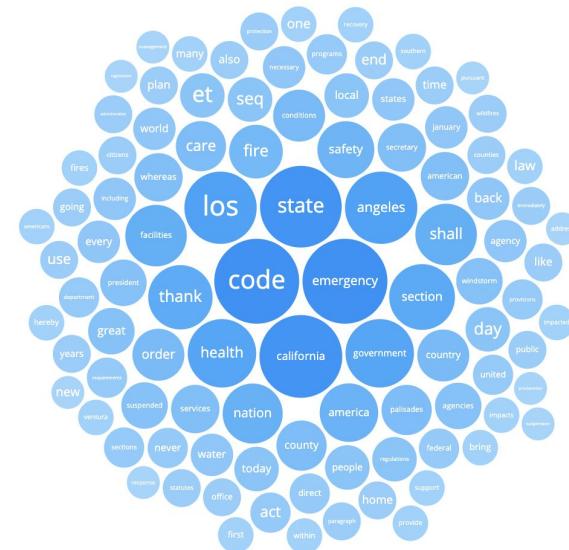
Lexos: Analyze > Top Words Example

Top Words		Generate		Download	
Document "Putting People over Fish" Compared To The Corpus		Document "The Inaugural Address" Compared To The Corpus		Document "January 13, 2025 Statement from President Joe Biden" Compared To The Corpus	
water	8.5233	code	-4.2133	across	7.6677
fish	5.8646	california	-4.1058	wildfires	6.0576
secretary	5.3489	state	-3.5239	support	4.9187
commerce	5.0774	emergency	-3.3722	firefighting	4.7673
interior	5.0774	section	-3.2581	suppress	4.7673
memorandum	5.0774	facilities	-3.0467	federal	4.3881

Lexos: Visualize



Word Cloud: visualize a word cloud across the entire text/corpus.



Bubbleviz: visualize word counts through bubbles across the entire text/corpus.

Lexos: Visualize > Multicloud



Voyant vs. Lexos: Word clouds

How does the Voyant word cloud below compare to the one made using Lexos?



Lexos Word cloud



What could be causing this distinction?

Lexos: Rolling Window

Rolling windows allow you to look at word trends across **one** document. To use a rolling window, first select a single text in the "Manage" screen, then:

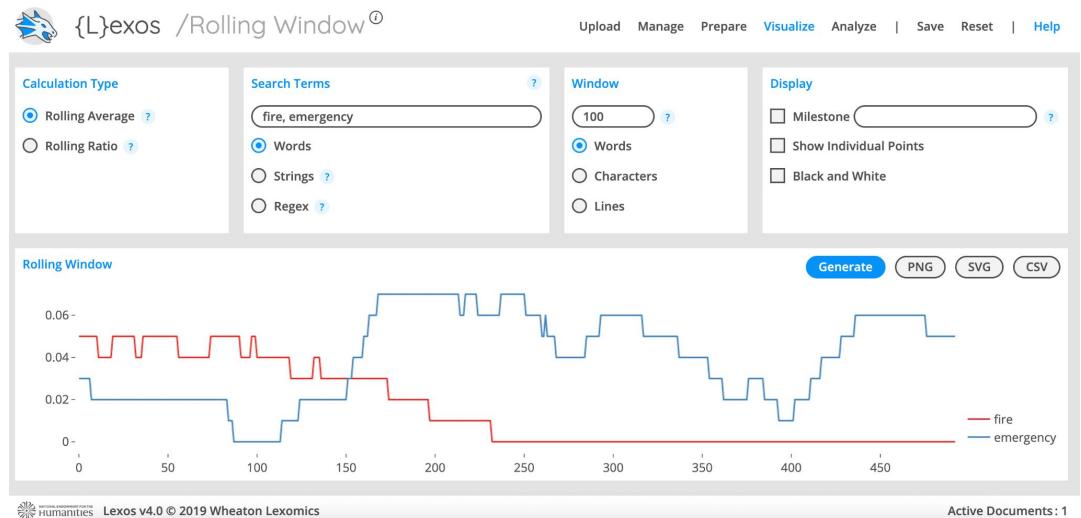
1. Go to “Visualize-> Rolling Window” and type in a search term you want to visualize. You can also search multiple terms by clicking “String” and separating words with a comma (climate, action)
2. Choose a Window size (the number of words each “window” contains). For shorter documents, it’s good to have a number like 300/500. For larger documents, you may want to make your window larger. Play around with the window size until you get a visualization that makes sense.
3. Click “Generate”

Lexos: Rolling Window Results

Using the document

“Proclamation of a State of Emergency”, and searching for the words ‘fire’ and ‘emergency’ with a window of 100, we can get an idea of how these terms work together in the document.

What other terms should we try?



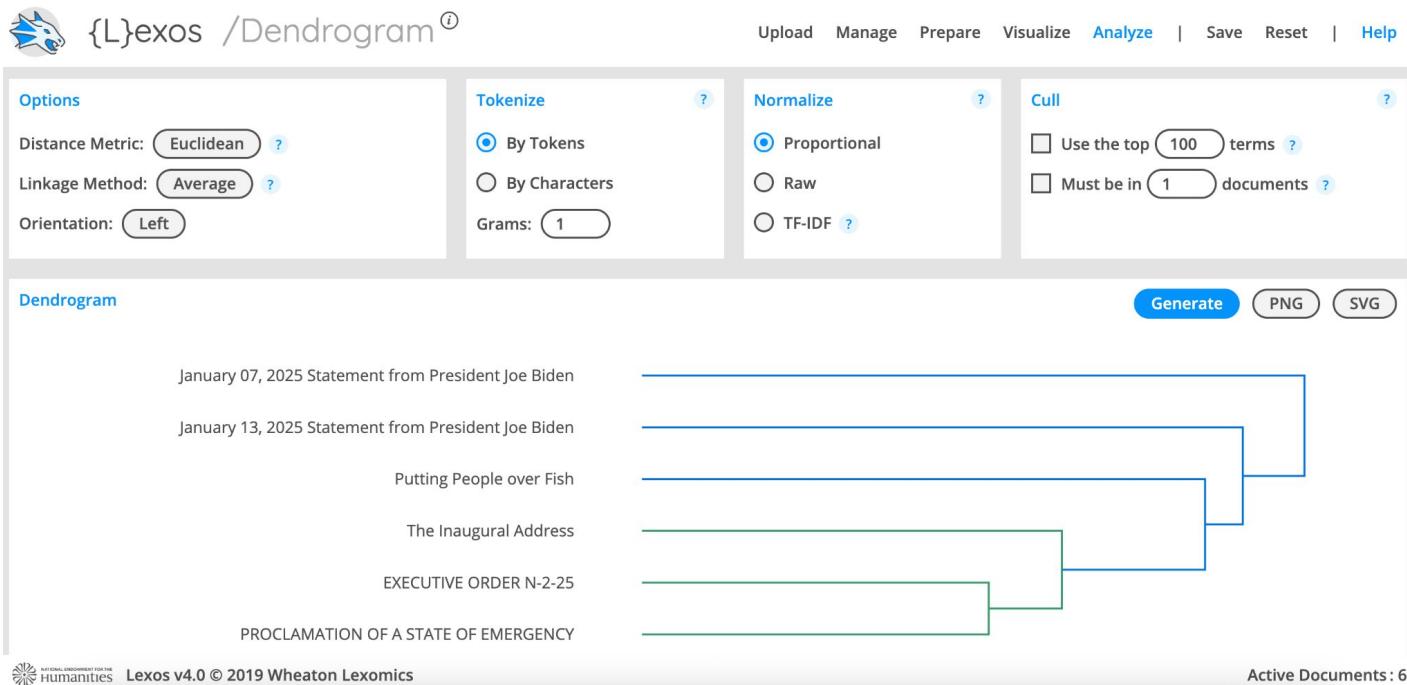
Lexos: Analyze > Dendrogram

The dendrogram demonstrates similarity between the different documents. Dendrograms require at least two documents to compare. Dendrograms are able to show the hierarchy between objects. Dendrograms show:

- Similarities between texts
 - The greater the distance between texts, the less similar they are
 - The smaller the distance between texts, the more similar they are

Lexos: Dendrogram

The dendrogram demonstrates similarity between the different documents.



Lexos: Save or Reset Your Results

Lexos allows you to **save** your results as a Lexos file. If you do this, you can re-upload the Lexos file any time to access your cleaned-up corpus as well as the different analyses you've done. You can also download modified text files from the “Manage” page—and you can even use those downloaded text files with other tools!

You can also save individual visualizations as images (PNGs).

Finally, if you want to start over, you can “Reset” your Lexos dashboard.

Your Turn!

Use the sample or other texts and begin practicing web-browser text analysis.

- Word Counter: <https://databasic.io/en/wordcounter/>
- Word Tree: <https://www.jasondavies.com/wordtree/>
- Voyant: <https://voyant-tools.org/>
- Lexos: <http://lexos.wheatoncollege.edu/upload>

Discussion Prompts

- What interesting or surprising results came up? What limitations are you observing?
- What kinds of texts would you be curious about comparing?
- Which features do you think will be useful in your future work?

A Brief Introduction to Web Scraping

Slide content courtesy of [Alyssa Smith](#)

*Feel free to ask questions at any point
during the presentation!*

Why Access Internet Data?

- Internet data can give us a way to (very imperfectly) quantify people's social lives online.
 - What are people talking about?
 - Who do people interact with?
 - How do communities form?
- It is especially useful at large scales.
 - Getting this kind of information on how people associate without social media data would be very difficult, if not impossible!
- Internet data is very rich in terms of context, content, and usability.
- Internet data captures certain times, cultures, and social contexts.
This is useful when researching recent and current issues.

How can you access internet data?

Unless you want to hand copy the contents of each web page, one at a time, you will need to use a program for automatically extracting data from the web. In some cases, websites provide their own tools, called **APIs**, that are designed to let you retrieve data that you specify. In other cases, you might use general software for **scraping** the contents of websites.

It helps to understand the general principles of how APIs and web scraping work, but typically each site will have its own specifications that you will need to learn to access their data.

Access Web Data Through APIs

- An API is a way for computer programs to talk to each other.
- APIs are code wrappers, a clean way to code communication with websites that eliminates the need for more complicated scraping.
- If you are trying to get a lot of information repeatedly from somebody else's computer program, an API is the way to do it!
- This might look like:
 - An analysis of all reddit posts mentioning "electric vehicles".
 - A program that emails you every time your elected officials in Congress post something with a negative sentiment.

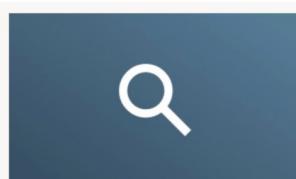
API Example - NY Times

- The New York Times features a Developer tool, available here: <https://developer.nytimes.com/>
- From here, users can sign up and access a variety of NYT content through their APIs.



Archive API

Get all NYT article metadata for a given month.



Article Search API

Search for New York Times articles.



Books API

Get NYT Best Sellers Lists and lookup book reviews.



Most Popular API

Popular articles on NYTimes.com.

[NYT Article Search API](#)

Web Scraping

- Sometimes websites don't have an API; you'll have to scrape the website.
- Scraping pulls the whole webpage—you then parse it and extract the data you want.
- This works better on structured websites that don't block bots (if you are scraping a website, you are a bot).
- Please obtain consent before scraping content from a site (or, at least, try to!)

Data privacy

- It's important to pay attention to data privacy when using digital resources
- At its simplest, **data privacy** is a person's ability to control what of their personal information is shared and with whom.
- To help you make informed decisions about interacting with digital tools in ways that honor your boundaries with your data and/or personal information, The DITI has prepared a handout on **Data Privacy**

Ethical Considerations of Scraping

- **Contextual Privacy**
 - When we think about privacy online we want to think of it as contextual. What someone might be comfortable saying in one context might not be something they would say to a researcher or want to be quoted in a publication.
- **Keeping People Safe**
 - It is risky to publicize the username, profile picture, or exact text of a social media post or profile.
 - To show example posts etc, you can make up your own or heavily redact them.

Learn More About Data Ethics & Text Mining

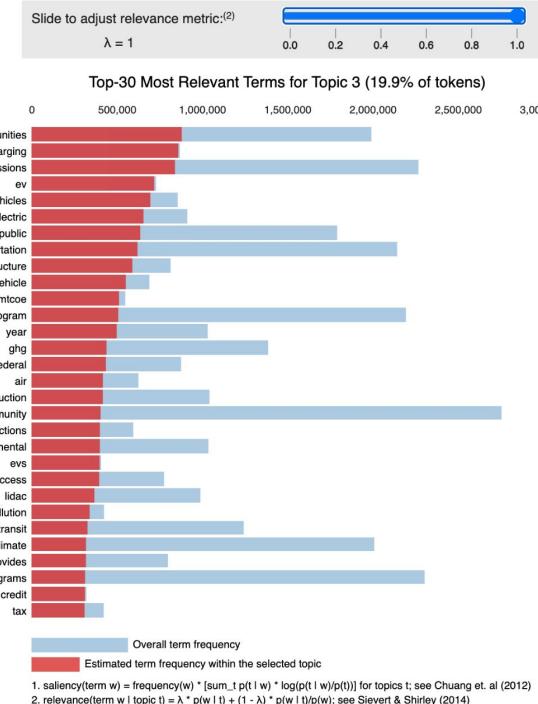
- [Data Ethics Handout](#)
- Northeastern Library [Guide on Text and Data Mining Library Databases](#)
- [ProQuest TDM](#)

Further Exploration

*Feel free to ask questions at any point
during the presentation!*

Further exploration: Topic Modeling

Topic modeling is a machine learning method that uses word co-occurrence within documents to identify "topics," or clusters of related terms. This is a topic model based on the Greater Boston Priority Climate Action Plan. In the visualization, topic 3 is selected.

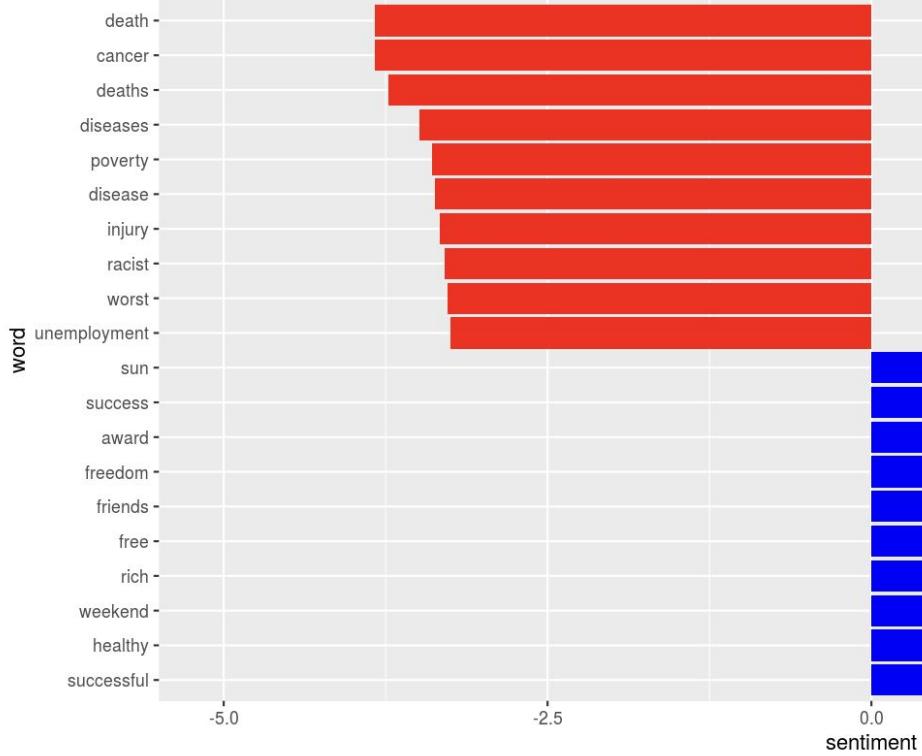


Topic model code generation assisted by ChatGPT and Gemini

Feel free to ask questions at any point during the presentation!

Further exploration sentiment analysis

Sentiment analysis uses dictionaries, and sometimes machine learning, to assign sentiment scores (e.g., positive and negative) to documents. You can try this out with the "[Drag and Drop Sentiment Analysis](#)" tool.



Greater Boston Priority Climate Action Plan

Feel free to ask questions at any point during the presentation!

For further exploration

- DITI handouts on [building a corpus](#) and more [links and resources](#) for text analysis
- NULab [list of resources for text analysis](#)
- [Programming Historian tutorials](#)
- [“Data-Sitters’ Club” tutorials](#)
- Library subject guides on text mining and analysis: [guide on getting started](#), [guide on vendor policies](#)

Thank you!

—Developed by Cara Marta Messina, Juniper Johnson, Sara Morrell, Ayah Aboelela, Jeff Sternberg, and Mel Williams

- For more information on DITI, please see: <https://bit.ly/diti-about>
- Schedule an appointment with us! <https://bit.ly/diti-meeting>
- If you have any questions, contact us at: nulab.info@gmail.com