

Computational Text Analysis for Content Analysis

Taught By Sara Morrell and Johan Arango-Quiroga
Digital Integration Teaching Initiative (DITI)

POLS 7346 Resilient Cities
Daniel Aldrich
Spring 2025

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, handouts, and data available at:

<https://bit.ly/sp25-aldrich-pols7346>

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:

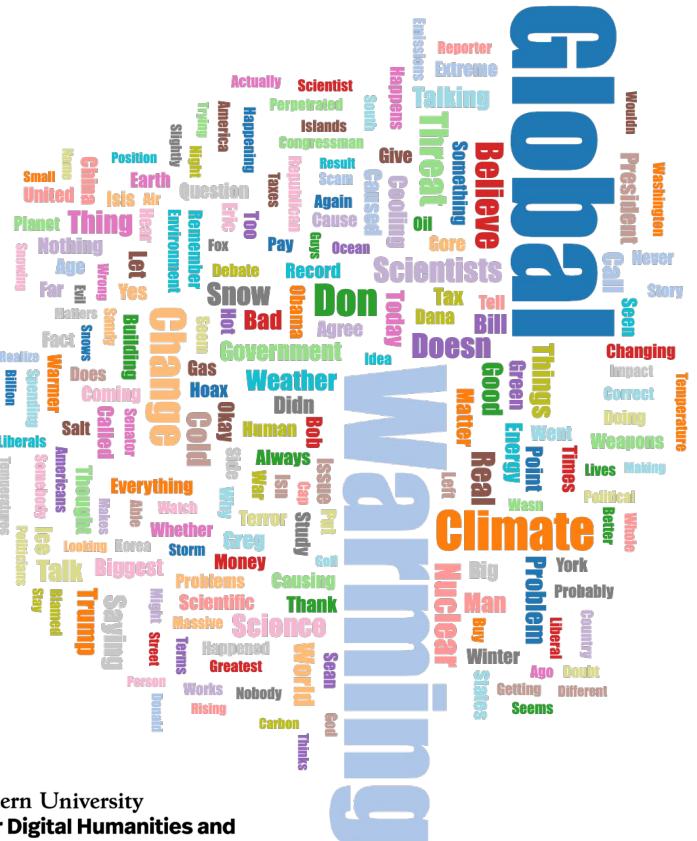
- Involves a computer drawing out patterns in a text, and a researcher interpreting those patterns.
- Includes methods such as word count frequency, keywords in context, computational modeling (with machine learning), and sentiment analysis.
- Is conducted using web-based tools or coding languages like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



Go to the [Television Explorer](#). Search “global warming,” “climate crisis,” “greenhouse effect.”

- What do you notice about the TV coverage of these terms over time?
What is surprising?
 - How do you think political values affects climate language?
 - 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.

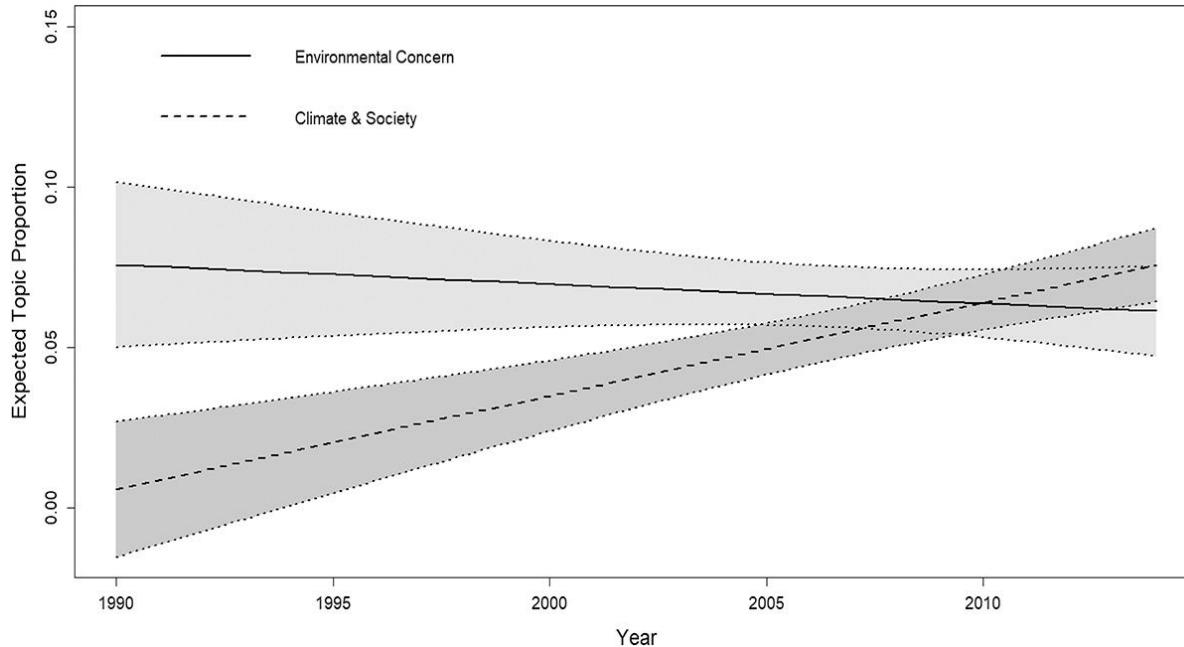
Key Terms (2/2)

- **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!*

Key Topics in environmental sociology



Topical prevalence of 'Environmental Concern' and 'Climate & Society' over time, 1990 – 2014 (with 95% confidence intervals).

Jeremiah Bohr & Riley E. Dunlap (2018) Key Topics in environmental sociology, 1990–2014: results from a computational text analysis, Environmental Sociology, 4:2, 181-195, DOI:

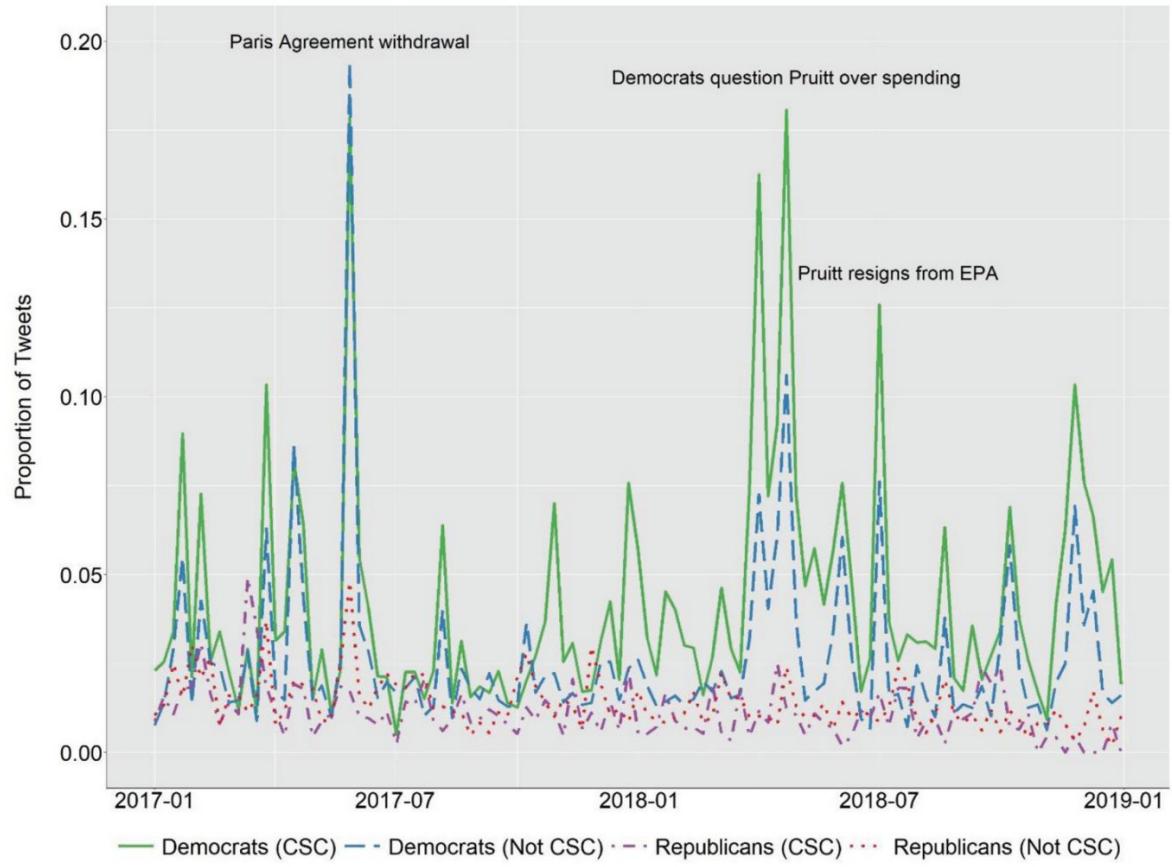
[10.1080/23251042.2017.1393863](https://doi.org/10.1080/23251042.2017.1393863)

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

U.S. Environmental Politics (1/2)

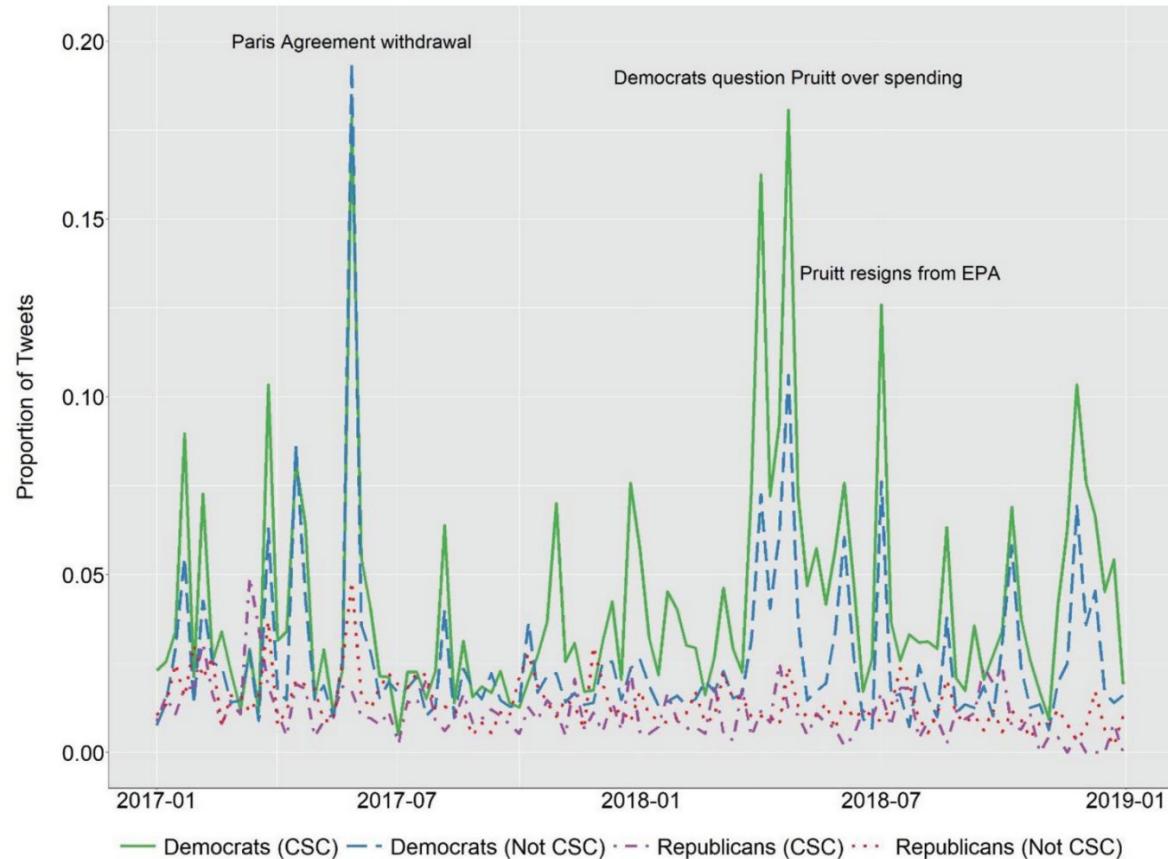
Weekly proportions of tweets discussing environmental issues sent by the 115th House of Representatives.

Key events and challenges:
a computational text
analysis of the 115th house
of representatives on
Twitter - Jeremiah Bohr in
Environmental Politics
(2021), 30 (3): 399-422

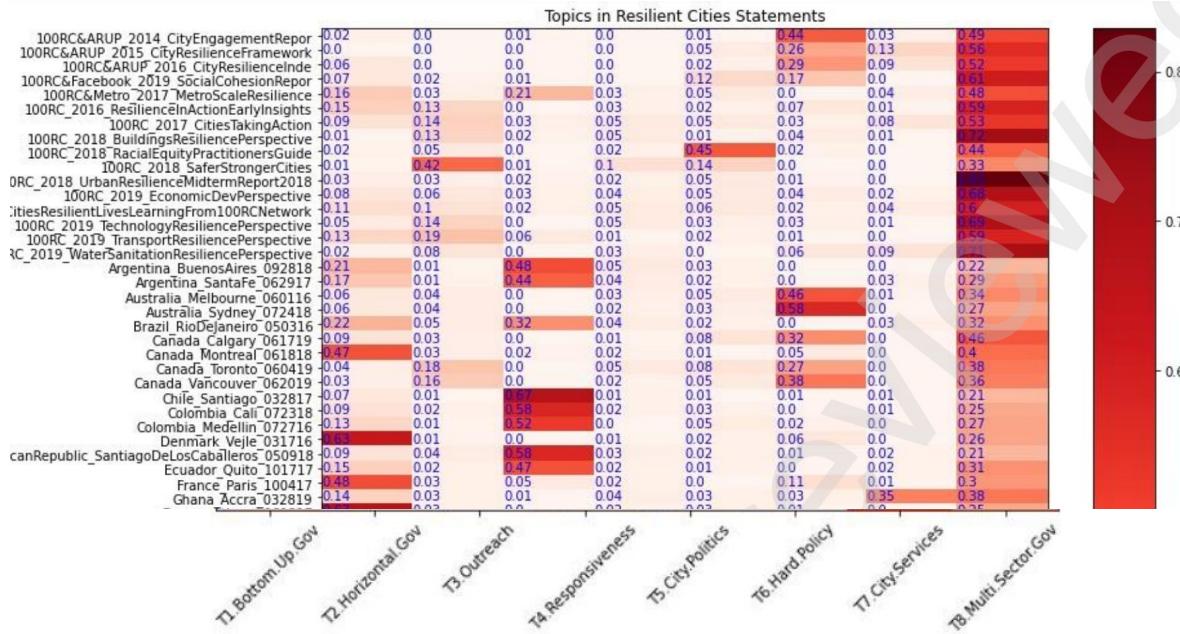


U.S. Environmental Politics (2/2)

To what extent do politicians publicly discuss environmental issues in line with public opinion and economic characteristics of their constituents?



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.

Boston Area Research Initiative (BARI) Data

You can use the [Boston Data Portal](#) and [BARI GIS map](#) to access Boston-specific information. This dataset includes 311 data, 911 calls, Airbnb listings, Craigslist housing ads, and more.

For example, BARI researchers looked at 911 calls and weather data [to identify heat islands and related illnesses around Boston](#), and consider how policies may be improved to reduce heat islands in Boston.

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 three metropolitan climate action plans published between 2023 and 2024:

- [Greater Boston Priority Climate Action Plan, March 2024](#)
- [PlaNYC Getting Sustainability Done, April 2023](#)
- [Priority Climate Action Plan for the Chicago Metropolitan Statistical Area, March 2024](#)

Sample Corpus

The sample .txt files are available on:

<https://bit.ly/sp25-aldrich-pols7346>

- For each file, click “Raw” in the top right corner.
- Right-click (PC) or Ctrl-click (Macs) on the text and choose “Save As.”
- Save as a .txt file on your computer.

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 **basic 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 Greater Boston Priority Climate Action Plan?



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 Boston’s climate action plan 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
energy	398	in the	147	advisory group members	54
emissions	213	of the	117	in the region	39
municipal	183	https www	103	the greater boston	28
community	182	the region	91	https www mass	26
transportation	182	for the	82	www mass gov	26
program	172	advisory group	81	department of energy	20
programs	168	new hampshire	75	ghg emissions reductions	19
communities	162	ghg emissions	75	measure advisory group	19
buildings	155	renewable energy	62	climate action plan	18

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 "Boston" from Boston's climate action plan

Emissions Reduction and Disclosure," The City of Boston, <https://www.boston.gov/departments/environment/building-emissions-reduction-and-disclosure> 32.
"PowerCorpsBOS," City of Boston, <https://www.boston.gov/departments/workforce-development/powercorpsbos>
33. "Good Jobs Metro Boston Coalition," City of Boston, <https://www.boston.gov/worker-empowerment/goodjobs-metro-boston-coalition> 34.
"Boston kicks off rollout of 750 electric school buses with first of the fleet transporting students," Electrek, 24 April, 2023, <https://electrek.co/2023/04/24/boston-rollout-750-electric-school-buses-first-fleet-transporting-students/>
35. "MetroWest Climate Equity Project," Metropolitan Area Planning Council, <https://www.mapc.org/resource-library/metro-west-climate-equity/> 36.
LIDAC is a term used by the Environmental Protection Agency (EPA). The EPA defines LIDAC as "any community that meets at least one of the following characteristics: 1) Identified as disadvantaged by the Climate and Economic

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.

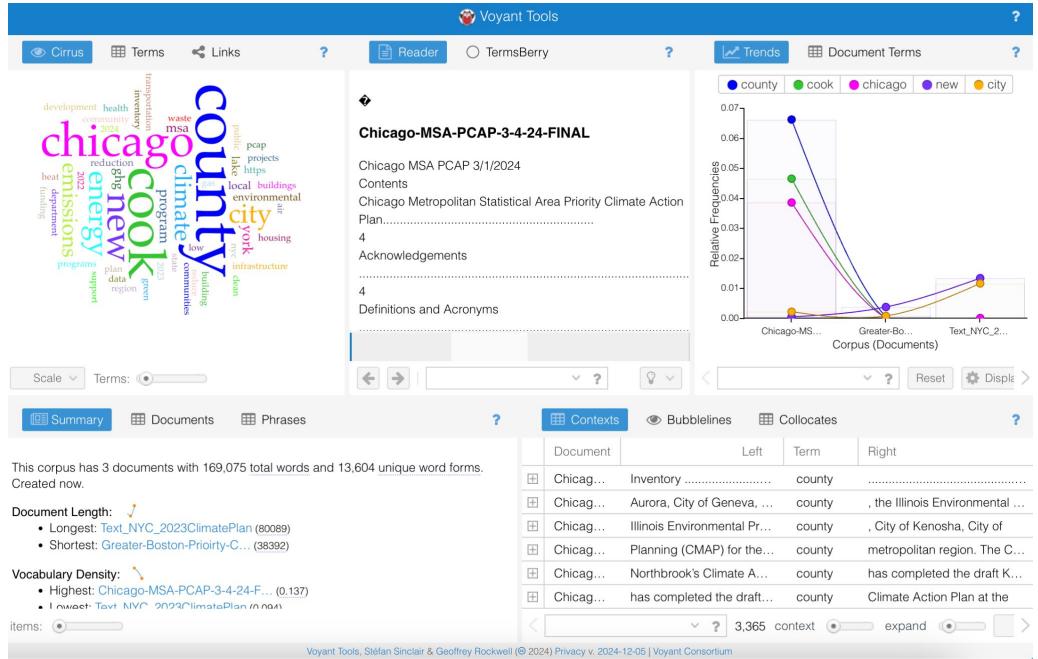
Click here for help and advanced options

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 two screenshots of the Voyant interface. On the left, the 'Terms' tab is active, displaying a word cloud with various terms like 'chicago', 'new', 'city', 'climate', etc., in different sizes. A red dashed arrow points from the top right corner of this pane towards the top right corner of the right-hand pane. On the right, the 'Collocates' tab is active, showing a table of collocation pairs and their counts:

Term	Collocate	Count (context)
chicago	chicago	12305
county	county	12167
cook	county	11102
county	cook	11004
cook	cook	8858
county	lake	1826
county	dupage	634
new	york	594
county	kane	544
new	city	428
city	new	426

Voyant: Contexts (concordances)

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

The screenshot shows the Voyant Contexts interface. At the top, there are three tabs: 'Contexts' (which is selected and highlighted in blue), 'Bubblelines', and 'Collocates'. To the right of the tabs is a help icon (a question mark). Below the tabs is a table with the following columns: 'Document', 'Left', 'Term ↓', and 'Right'. The table lists nine rows, each corresponding to a document named 'Text_NYC...'. The 'Left' column contains various names and titles, and the 'Right' column contains the word 'climate'. The 'Term ↓' column is centered above the table. At the bottom of the interface, there is a search bar containing the word 'climate' with a red rectangular border around it. To the right of the search bar are buttons for '872 context', 'expand', 'Scale', and a dropdown menu.

Document	Left	Term ↓	Right
Text_NYC...	support by Thornton Tomasetti nyc.gov/	climate	
Text_NYC...	Staff, Office of the Chief	climate	Officer, Mayor's Office ...
Text_NYC...	Kathy Hochul Peter Davidson, Aligned	climate	Capital Santos Rodrigu...
Text_NYC...	Amy Turner, Sabin Center for	climate	Change Law Arif Ullah,...
Text_NYC...	Falade, Meena Ardebili, Sarah Varghese	climate	CABINET MEMBERS A...
Text_NYC...	EXECUTIVE DIRECTOR, MAYOR'S OFFICE OF	climate	& ENVIRONMENTAL J...
Text_NYC...	Public Safety, Phil Banks Chief	climate	Officer, Rohit T. Aggarw...
Text_NYC...	Bay – Rockaway Parks Conservancy. (n.d.).	climate	and Community Resilie...
Text_NYC...	York City Mayor's Office of	climate	& Environmental Justic...

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

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

Voyant: Your Turn!

Use the sample text or texts of your choice and begin practicing web-browser text analysis. **Explore Voyant features!**

Discussion Prompts

- What interesting or surprising results came up?
- How might you interpret those results based on what you know about current climate plans?
- What other kinds of documents would be useful to compare in Voyant to research trends in climate and disaster planning?

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.

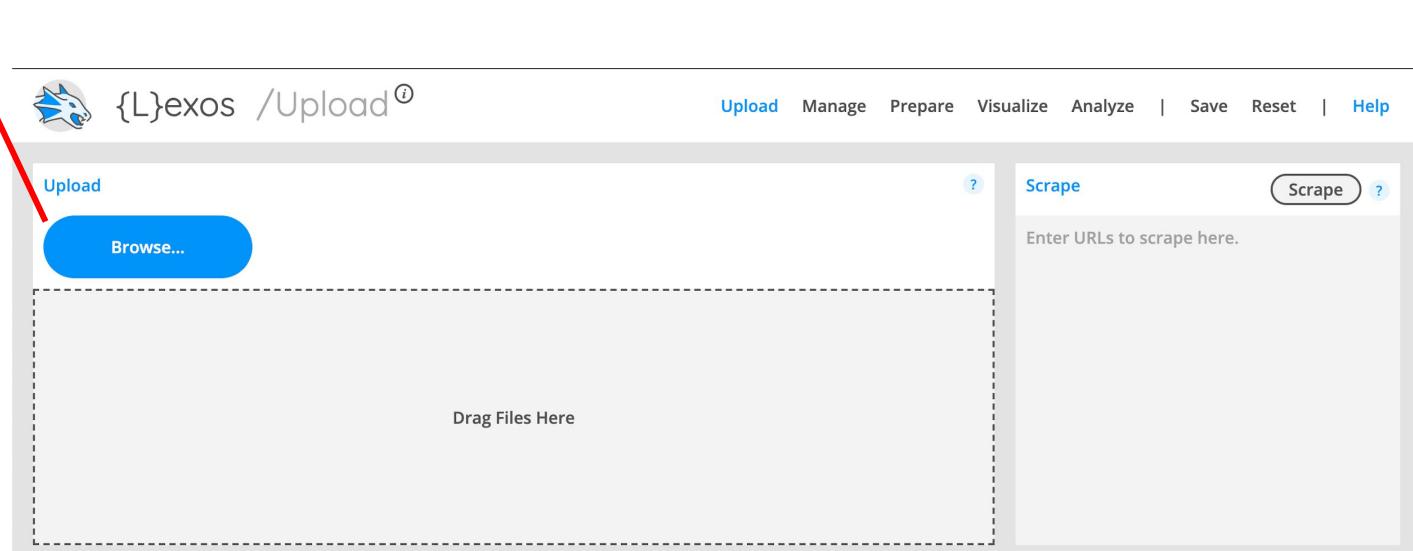
- **Upload:** upload your .txt file
- **Manage:** select the files you want to prepare and analyze
- **Prepare:** prepare your text for analysis
- **Visualize:** create visualizations of patterns across your corpus or in single texts
- **Analyze:** analyze your text

<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)

{L}exos /Manage^①

Upload **Manage** Prepare Visualize Analyze | Save Reset | Help

Active	#	Document	Class	Source	Excerpt
●	1	Greater-Boston-Priority-Climate-Action-Plan-March-2024		Greater-Boston-Priority-Climate-Action-Plan-March-2024.txt	Greater Boston Priority Climate Action Plan March 2024 Greater Boston Priority Climate Action Plan March 2024 PREPARED BY: M... ...ould be involved and that the region would need to see further investments in the MBTA to make this measure successful. 152
●	2	Chicago-MSA-PCAP-3-4-24-FINAL		Chicago-MSA-PCAP-3-4-24-FINAL.txt	Chicago MSA PCAP 3/1/2024 Contents Chicago Metropolitan Statistical Area Priority Climate Action Plan.....dith Makra, Director of Environmental Initiatives emakra@mayorscaucus.org 630-327-4193 184 Chicago MSA PCAP 3/1/2024 185
●	3	Text_NYC_2023ClimatePlan		Text_NYC_2023ClimatePlan.txt	The City of New York Mayor Eric Adams April 2023 Getting Sustainability Done PlaNYC Letter from the Mayor Introduction Our Visi... ...ncy analysis by WSP USA Graphical analysis by WXY Design by Nowhere Office Policy support by Thornton Tomasetti nyc.gov/climate

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 buttons for Upload, Manage, Prepare (which is highlighted in blue), Visualize, Analyze, Save, Reset, and Help. Below these are tabs for Previews, Preview, Apply, and Download. The main area has several sections: 'Scrubbing Options' with checkboxes for Make Lowercase (checked), Remove Digits (checked), Remove Spaces, Remove Tabs, Remove Newlines, and options for Scrub Tags, Remove Punctuation, Keep Hyphens, Keep Apostrophes, and Keep Ampersands; 'Lemmas' with a text input field and an Upload button; 'Consolidations' with a text input field and an Upload button; 'Stop/Keep Words' where the 'Stop' radio button is selected (highlighted with a red box); a list of stop words including cook, county, i, me, my, myself, we, our, ours, ourselves; 'Special Characters' with radio buttons for None (selected), MUFI 3, Early English HTML, MUFI 4, Old English SGML, and a text input field for entering special characters. On the right side, there are preview sections for 'Greater-Boston-Priority-Climate-Action-Plan-March-2024' and 'Chicago-MSA-PCAP-3-4-24-FINAL', and a 'Text_NYC_2023ClimatePlan' section. At the bottom left is the National Endowment for the Humanities logo and 'Lexos v4.0 © 2019 Wheaton Lexomics'. At the bottom right is 'Active Documents: 3'.

You can also add stop words particular to your corpus.

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

Lexos: Applying your Preparations

BEFORE PREP

Previews

Preview

Apply

Download

[Greater-Boston-Priority-Climate-Action-Plan-March-2024](#)

Greater Boston Priority Climate Action Plan March 2024

Greater Boston Priority Climate Action Plan March 2024

PREPARED BY: M... ...ould be involved and that the region would need to see further investments in the MBTA to make this measure successful. 152

AFTER PREP

Previews

Preview

Apply

Download

[Greater-Boston-Priority-Climate-Action-Plan-March-2024](#)

greater boston priority climate action plan march greater boston priority climate action plan march prepared metropolitan... ...standard commuter hours said riders union employers involved region would need see investments mbta make measure successful

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 "Greater-Boston-Prioirty-Climate-Action-Plan-March-2024" Compared To The Corpus

chicago -21.7793

city -13.9094

municipal 13.5272

massachusetts 12.8756

na -11.712

Document "Chicago-MSA-PCAP-3-4-24-FINAL" Compared To The Corpus

chicago 43.5797

lake 21.0948

msa 17.0505

new -16.9081

pcap 14.1964

Document "Text_NYC_2023ClimatePlan" Compared To The Corpus

chicago -30.8436

msa -15.1287

lake -14.9818

new 14.0815

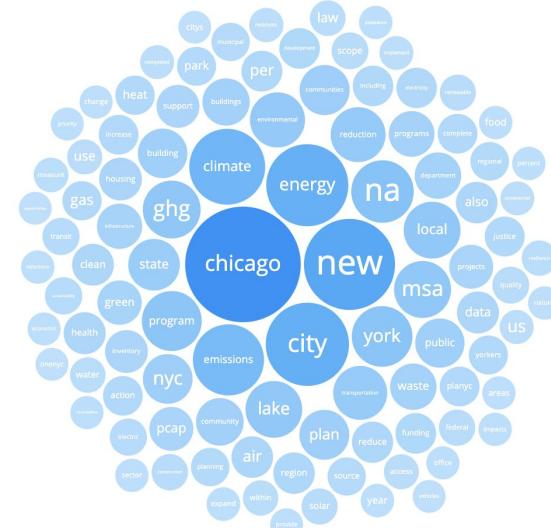
city 13.0095

Lexos: Visualize

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



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



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

Lexos: Visualize > Multicloud

{L}exos /Multicloud^①

Upload Manage Prepare **Visualize** Analyze | Save Reset | Help

Multicloud

Font: Open Sans Term Count: 100 Color: Default Generate

Greater-Boston-Priority-Climate-Action

PNG SVG

Chicago-MSA-PCAP-3-4-24-FINAL

PNG SVG

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?

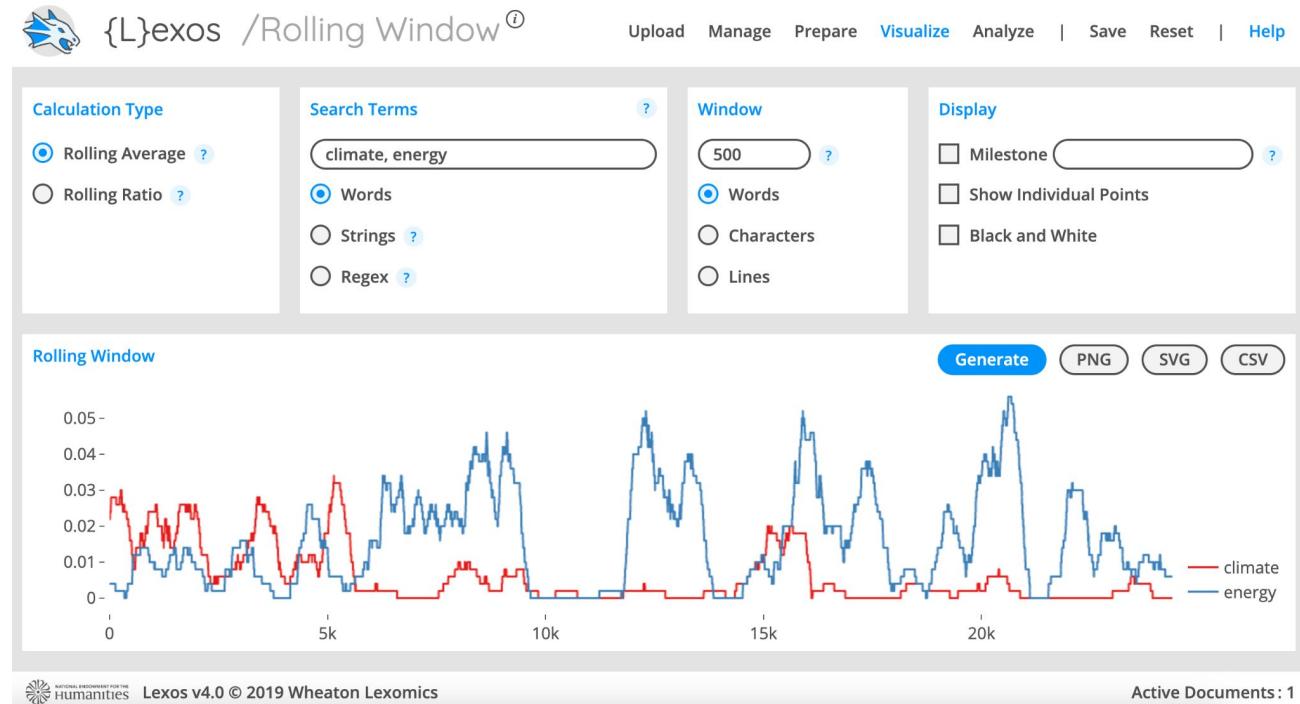
Lexos: Rolling Window

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 Boston's climate action plan, and searching for the words 'climate' and 'energy' with a window of 500, we can get an idea of how these terms work together in the document.



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

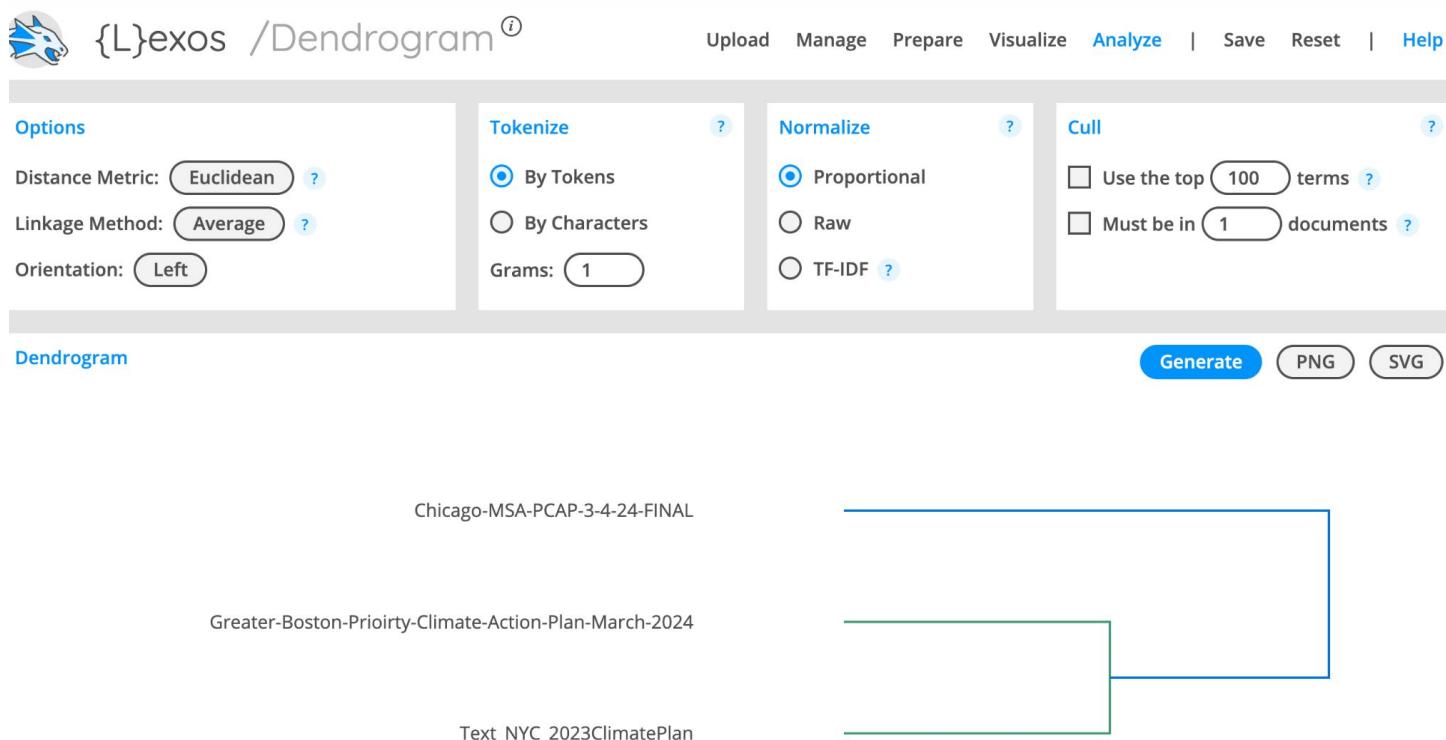
Lexos: Analyze > Dendrogram

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.



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

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.

Lexos: Your Turn!

Use the sample text or texts of your choice and begin practicing web-browser text analysis. **Explore Lexos's features!**

Discussion Prompts

- Between Voyant and Lexos, which tool do you prefer and why?
- How would you want to use these tools for your research? Which features do you think will be useful in your analysis?
- How do you think this computational text analysis can complement your other research methods?

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.

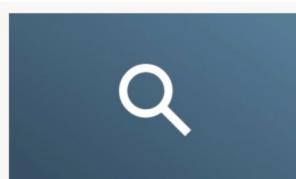
API Example - NY Times

- 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 Web Scraping

- <https://bit.ly/ScrapingSlides>
- [Data Ethics Handout](#)
- Northeastern Library [Guide on Text and Data Mining Library Databases](#)
- Databases for learning and applying text analysis:
 - [Constellate](#)
 - [ProQuest TDM](#)

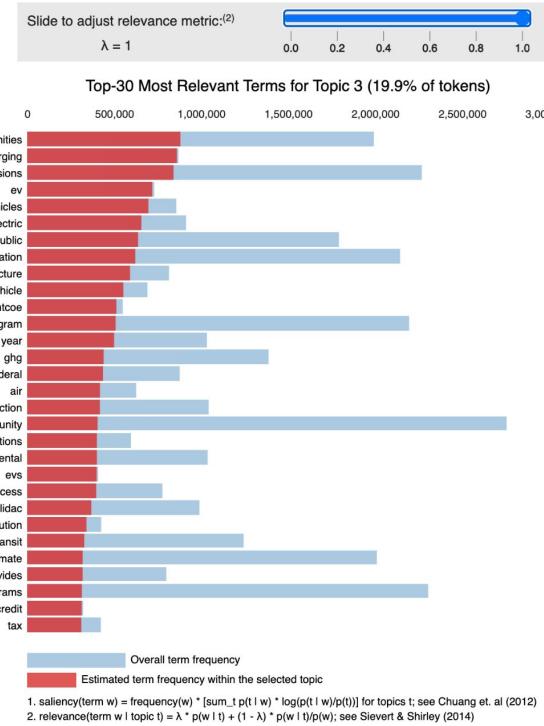
Further Exploration

*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.

Selected Topic: 3 Previous Topic Next Topic Clear Topic

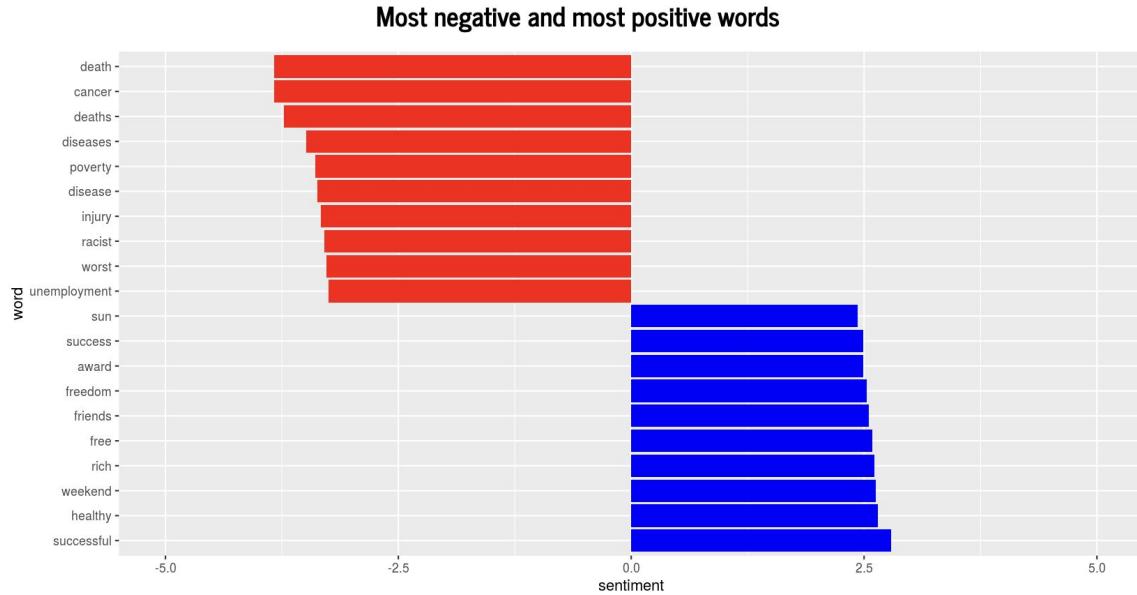


Feel free to ask questions at any point during the presentation!



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, and Jeff Sternberg

- 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