



No Fear Digital Humanities!

Text Mining Historical Documents with the Gale Digital Scholar Lab

10 March 2022

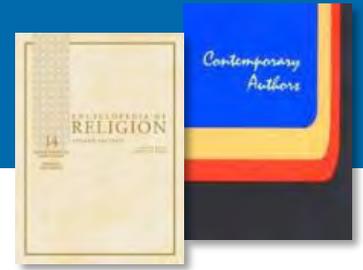
Masaki Morisawa
Senior Product Manager, Gale Asia



About Gale

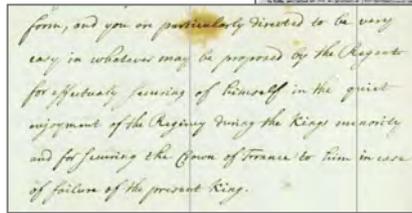
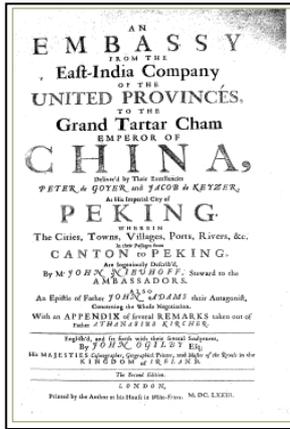


- HQ in Farmington Hills, Michigan, USA
- One of the foremost library reference publishers
- Has many well-established imprints, including:
 - Gale
 - Charles Scribner's Sons
 - Macmillan Reference USA
- Publishing formats include:
 - **Print** library reference, such as thematic encyclopedias, annual directories, literary biographies/criticisms, etc.
 - **eBook** versions of print publications and an eBook platform
 - Aggregated **Journal databases**
 - **Subject-specific Databases** combining various content together
 - **Microform** collections and serials
 - **Gale Primary Sources**: digital archives of historical material



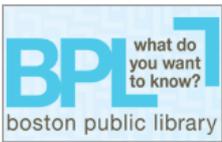
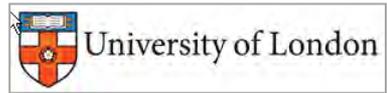
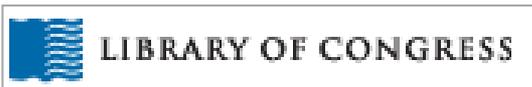
What are Gale Primary Sources?

Gale Primary Sources



- Gale has spent the past 45 years building one of the world's largest scholarly primary source online libraries, spanning 560 years of global history
- More than 200 digital collections, comprising more than 179 million pages, including:
 - Nineteenth Century Collections Online
 - The Times Digital Archive
 - China and the Modern World
 - U.S. Declassified Documents Online
 - Archives of Sexuality and Gender
 - *And many more!*

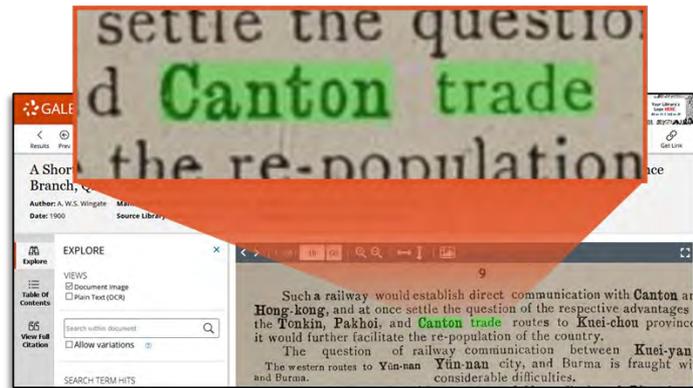
Examples of Original Holding Institutions



We scan and digitize valuable historical works



Enable the discovery of hidden words

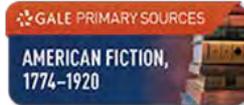
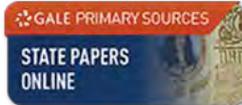
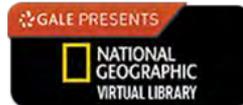


OCR = Optical Character Recognition

HTR = Handwritten Text Recognition*

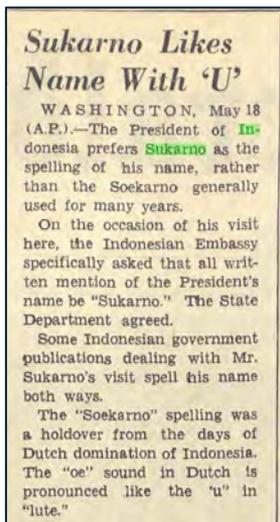
* HTR available in parts of *China and the Modern World* only

We have *many* such products



What is the Gale Digital Scholar Lab?

Some Researchers are Interested in the OCR Text behind the Images



Page Images (OCR not seen)

Sukarno Likes Name with 'U' Sukarno Likes Name With 'U', WASHINGTON, May 18 (A.P.).—The President of Indonesia prefers Sukarno as the spelling of his name, rather than the Soekarno generally used for many years. On the occasion of his visit here, the Indonesian Embassy specifically asked that all written mention of the President's name be "Sukarno." The State Department agreed. Some Indonesian government publications dealing with Mr. Sukarno's visit spell his name both ways. The "Soekarno" spelling was a holdover from the days of Dutch domination of Indonesia. The "oe" sound in Dutch is pronounced like the "u" in "lute."

OCR Text (generated from image)

What If you could combine hundreds of text files like these...

Sukarno Likes Name with 'U' Sukarno Likes Name With 'U' WASHINGTON, May 18 (A.P.)—The President of Indonesia prefers Sukarno as the spelling of his name, rather than the Soekarno generally used for many years. On the occasion of his visit here, the Indonesian Embassy specifically asked that all written mention of the President's name be 'Sukarno.' His State Department agreed. Some Indonesian government publications dealing with Mr. Sukarno's visit spell his name both ways. The 'Soekarno' spelling was a holdover from the days of Dutch domination of Indonesia. The 'oe' sound in Dutch is pronounced like the 'ou' in 'lute.'



Eisenhower, Sukarno Agree on 'Freedom' 881 mm President Sukarno of Indonesia, and Mrs. Richard M. Nixon at a Washington reception. Eisenhower, Sukarno see on 'Freedom' Ag. By the United Press WASHINGTON, May 20.—President Eisenhower said Friday night that the people of the United States must never forget that unless we defend the freedom of the other fellow we are endangering our own. Mr. Eisenhower, in a toast to visiting Indonesian President Sukarno, said that he wanted to express his gratitude to Mr. Sukarno for recalling a basic American responsibility. 'I want to thank you for what you have reminded us of.'



Sukarno Hailed on Broadway Sukarno Hailed On Broadway NEW YORK, May 23 (A.P.)—New York today accorded President Sukarno of Indonesia the traditional honor for visiting dignitaries—a ticker tape parade along Lower Broadway through the skyscrapers of the financial district. The President stood in the rear of an open car for the 15-minute drive from Battery Park to City Hall. A light rain began to fall just before the parade started, but even so police estimated the crowd along the line of march and in City Hall Plaza at 80,000. Many persons along the parade route waved the red and white flag of Indonesia. Mr. Sukarno waved repeatedly. Arriving at City Hall, Mr. Sukarno was greeted by Mayor Robert F. Wagner. Both men stood at attention as the national anthems of the United States and Indonesia were played.



Sukarno and Pope Hold Talk Sukarno And Pope Hold Talk By Barrett McGurn From the Herald Tribune Bureau g) 1956 New York Herald Tribune, Inc., ROME, June 13.—Pope Plus XII promised President Achmed Sukarno of Indonesia here today that foreign missionaries of the Roman Catholic Church will be withdrawn from his country as quickly as possible, and appealed in effect for 'patience' during what he calls the 'transition period.' The Indonesian revolutionary nationalist leader assured the Pope in reply that he considered the 1,000,000 Indonesian Catholics, slightly more than 1 per cent of the population, a loyal part of his country and wished the Catholic religion well in his homeland.



...



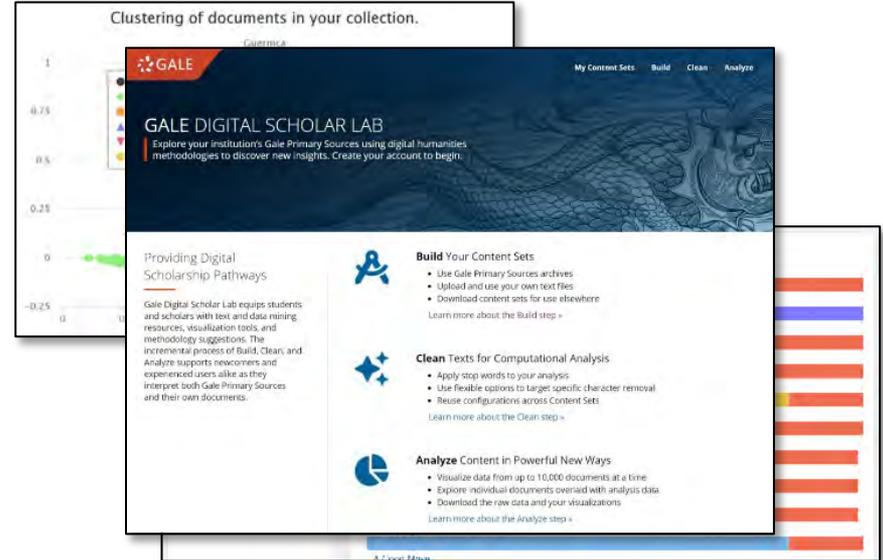
Bag of Words



... Require Different Platforms



Gale Primary Sources

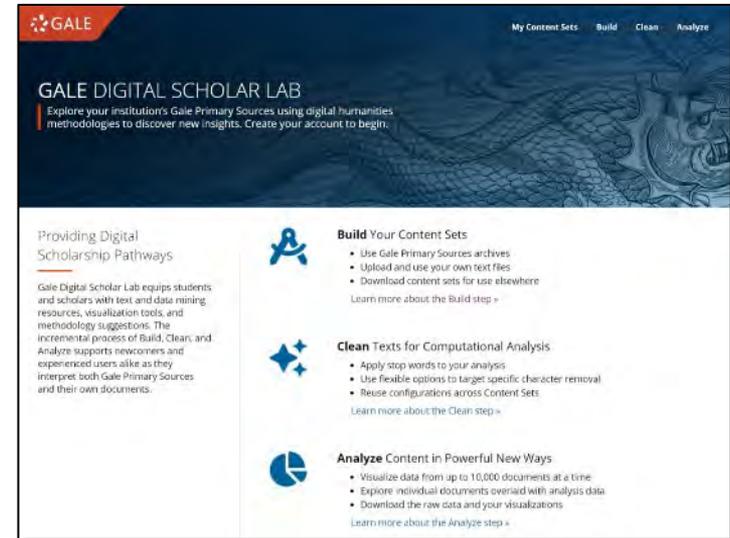


Gale Digital Scholar Lab

Gale Digital Scholar Lab

Gale's Cloud-based Text and Data Mining Research Environment

- **Access** OCR/HTR text from Gale Primary Sources
- **Build** custom content subsets
- **Clean** OCR text to reduce or control noise
- **Analyze** with simple but powerful tools
- **Export** statistical data, visualizations, OCR text in standard formats
- **Organize** and manage research



The screenshot displays the Gale Digital Scholar Lab interface. At the top, there is a navigation bar with the GALE logo and links for 'My Content Sets', 'Build', 'Clean', and 'Analyze'. Below the navigation bar, the main heading reads 'GALE DIGITAL SCHOLAR LAB' with a sub-heading: 'Explore your institution's Gale Primary Sources using digital humanities methodologies to discover new insights. Create your account to begin.' The interface is divided into three main sections, each with an icon and a brief description:

- Providing Digital Scholarship Pathways**: Accompanied by a book icon, this section describes how the lab equips students and scholars with text and data mining resources, visualization tools, and methodology suggestions. It notes that the incremental process of Build, Clean, and Analyze supports newcomers and experienced users alike as they interpret both Gale Primary Sources and their own documents.
- Build Your Content Sets**: Accompanied by a gear icon, this section lists three steps: 'Use Gale Primary Sources archives', 'Upload and use your own text files', and 'Download content sets for use elsewhere'. A link to 'Learn more about the Build step' is provided.
- Clean Texts for Computational Analysis**: Accompanied by a star icon, this section lists three steps: 'Apply stop words to your analysis', 'Use flexible options to target specific character removal', and 'Reuse configurations across Content Sets'. A link to 'Learn more about the Clean step' is provided.
- Analyze Content in Powerful New Ways**: Accompanied by a pie chart icon, this section lists three steps: 'Visualize data from up to 10,000 documents at a time', 'Explore individual documents overlaid with analysis data', and 'Download the raw data and your visualizations'. A link to 'Learn more about the Analyze step' is provided.

Gale Digital Scholar Lab Workflow and Tools

Gale Digital Scholar Lab: The Workflow



The Six Text Mining Tools Available in the Lab

Name of Tool	Common Uses	Library / Toolkit	Open Source?
Ngram	Extracts frequently occurring words and phrases	Lucene	No
Clustering	Groups documents in “clusters” according to their content	SciKit Learn	Yes
Topic Modelling	Extracts common topics found across many documents	Mallet	Yes
Sentiment Analysis	Classifies documents between positive to negative sentiment	OpenNLP	Yes
Named Entity Recognition	Extracts “Named Entities” (proper names, dates, prices, etc.)	spaCy	Yes
Parts-of-Speech Tagger	Compares the use of grammatical properties (nouns, verbs, etc.)	spaCy	Yes

1. Ngram (Frequent Terms)

The Ngram tool breaks down the text in your content set into sets of n consecutive words (Ngrams) and counts their occurrences. Ngrams are like *phrases* except that the resulting pieces are not always meaningful. They are helpful in finding **frequently occurring words and phrases** in your content set.

For example, if we split the text **“Nobody in Singapore drinks Singapore Slings.”**:

N=1: Unigram

“Nobody”, “in”, “Singapore”, “drinks“, ... etc.

N=2: Bigram

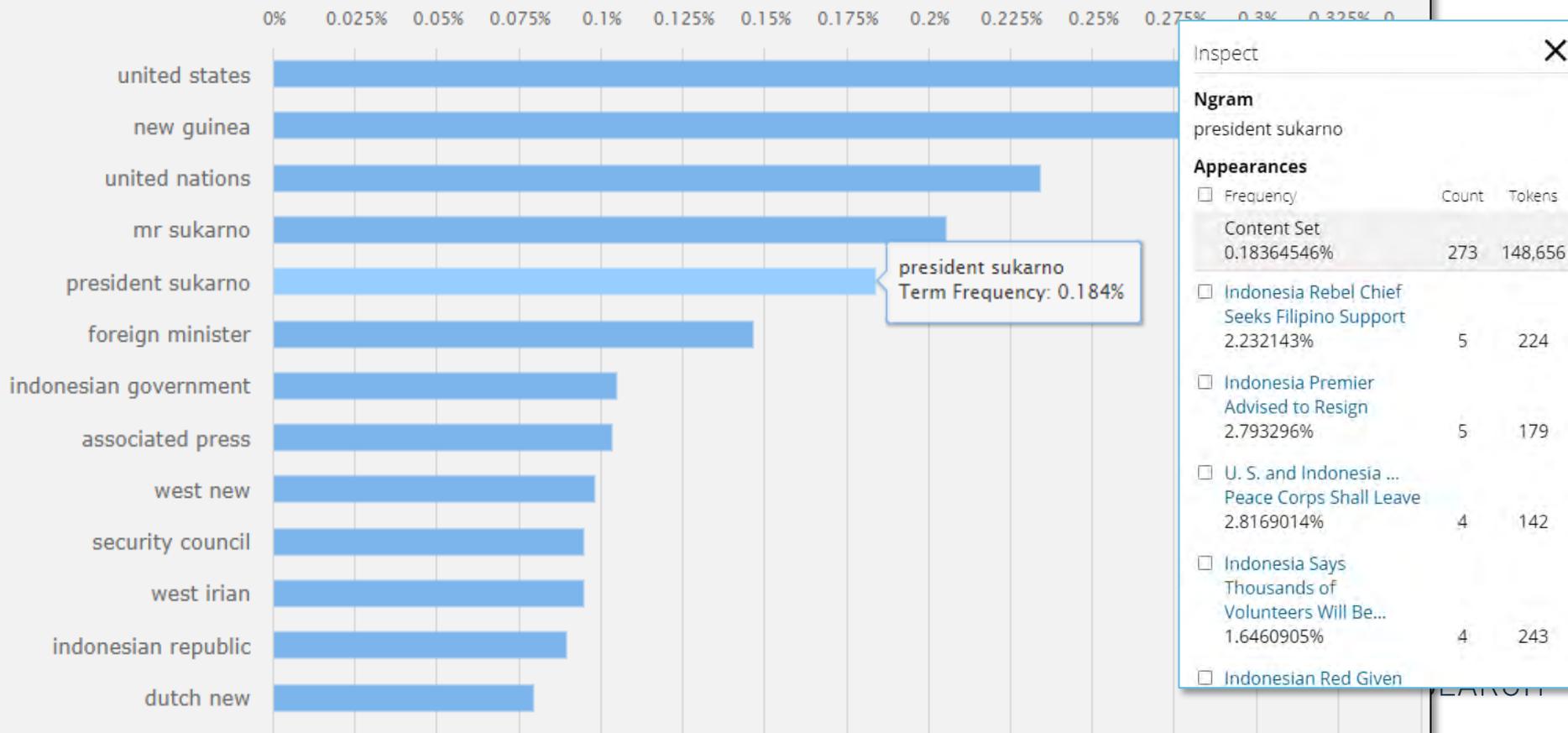
“Nobody in”, “in Singapore”, “Singapore drinks“, “drinks Singapore“, ... etc.

N=3: Trigram

“Nobody in Singapore”, “in Singapore drinks“, “Singapore drinks Singapore“, “drinks Singapore Slings“, ... etc.

Note: Ngrams that are composed entirely of stop words will not be considered. Ex. **“To Be or Not to Be”**

1. Ngram Lab Output (Bar Graph)



2. Named Entity Recognition (NER)

- NER recognises and extracts “Named Entities” (proper names, dates, prices, etc.) from documents within a content set, and output lists of such entities.
- Not designed to analyze textual features, but to extract them, and create new datasets which can then be used to contextualise content sets.
- Example “entity types”:
 - **people** (including fictional),
 - **groups** (nationalities, religious, or political),
 - **organizations** (companies, agencies, institutions),
 - **locations** (countries, states, cities),
 - **products** (objects, vehicles, foods, etc.), works of art (titles of books, songs, etc.),
 - **dates** (absolute or relative dates or periods)

2. Named Entity Recognition Output

Legend

View

Top 200 Entities by Count

Entity Search

Search for entities

Entity categories

Category

Date

Time

Geography

Geo-Political Entity

Place

Artwork

Event

Law

Product

Person

Entity ↕	Category ↕	Documents ↕	Count ↕
Japan	Geo-Political Entity	215	3396
Japanese	Cultural Group	206	2544
U.S.	Geo-Political Entity	129	1534
US	Geo-Political Entity	159	1510
1	Number	236	1155
2	Number	269	981
Okinawa	Organization	186	963
the United States	Geo-Political Entity	119	830
4	Number	214	745
3	Number	209	732
NSC	Organization	43	482
5	Number	165	458
Korea	Geo-Political Entity	75	417
China	Geo-Political Entity	82	397

List of Named Entity Categories

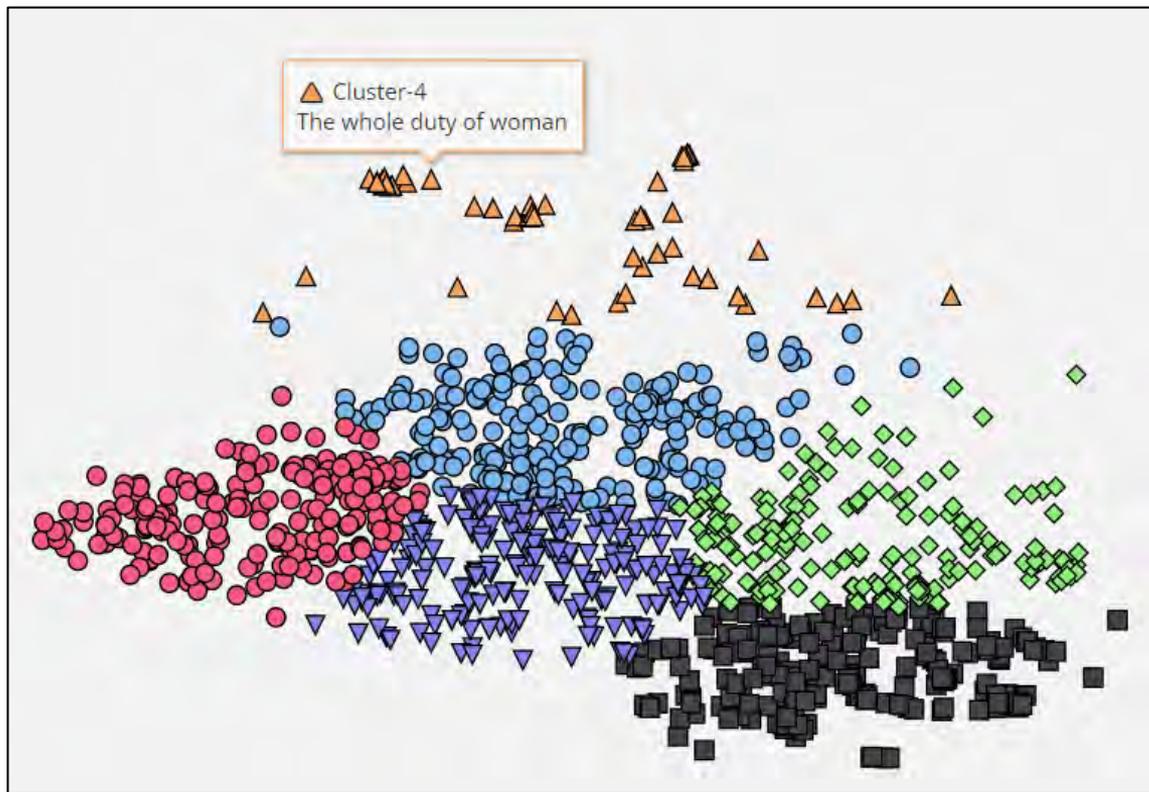
List of Named Entities found

3. Clustering

- The Clustering tool looks at a collection of documents, and groups similar documents into clusters
- The user specifies the number of clusters to display
- The tool analyses the documents using statistical methods to group them around particular features



3. Clustering Output



Inspect



Document

The polite instructor; or, youth's museum.
◆ Consisting of moral essays, tales, fables, visions, and allegories. Selected from the...

Open Document

Top 20 Most Similar Documents

The smaller the distance, the greater the similarity to the selected document. [Download Full List](#)

<input type="checkbox"/>	Distance	Title
<input type="checkbox"/>	0.0001	◆ Instructions for a young lady, in every sphere and period of life. Containing, I. A...
<input type="checkbox"/>	0.0001	◆ The lady's miscellany; or Pleasing essays, poems, stories, and examples, for...
<input type="checkbox"/>	0.0005	◆ The ladies complete letter-writer; teaching the art of inditing letters On every...
<input type="checkbox"/>	0.0013	◆ Instructions for a young lady,

4. Sentiment Analysis

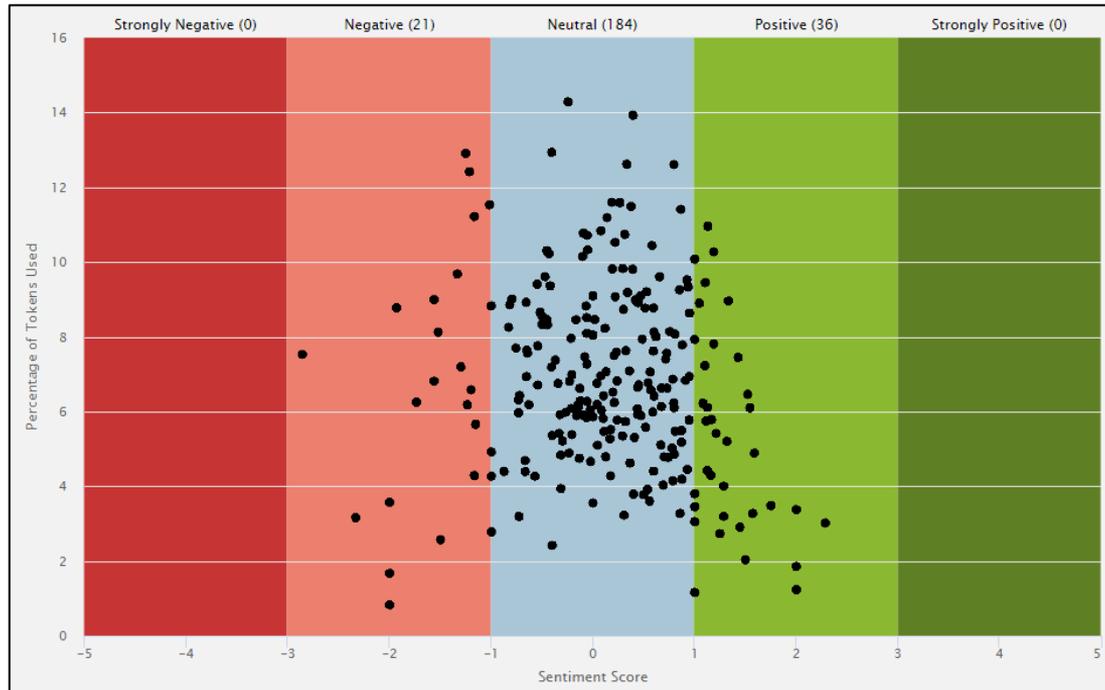
This tool calculates the “sentiment score” of each document, which indicates the use of positive or negative words within them

It works by applying a lexicon of terms and their associated sentiment scores to a piece of text

It also divides each document into 20 sections and calculates the sentiment score for each section



4. Sentiment Analysis Output



Inspect ✕

Document
-0.169 The politick wife: or, the devil outwitted by a woman

[Open Document](#) [Add to Content Set](#)

Document Subsections

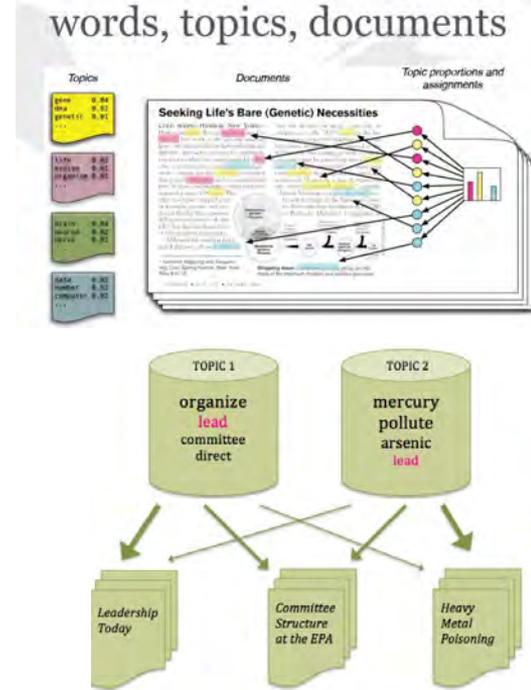
Beginning End

Percent Scored

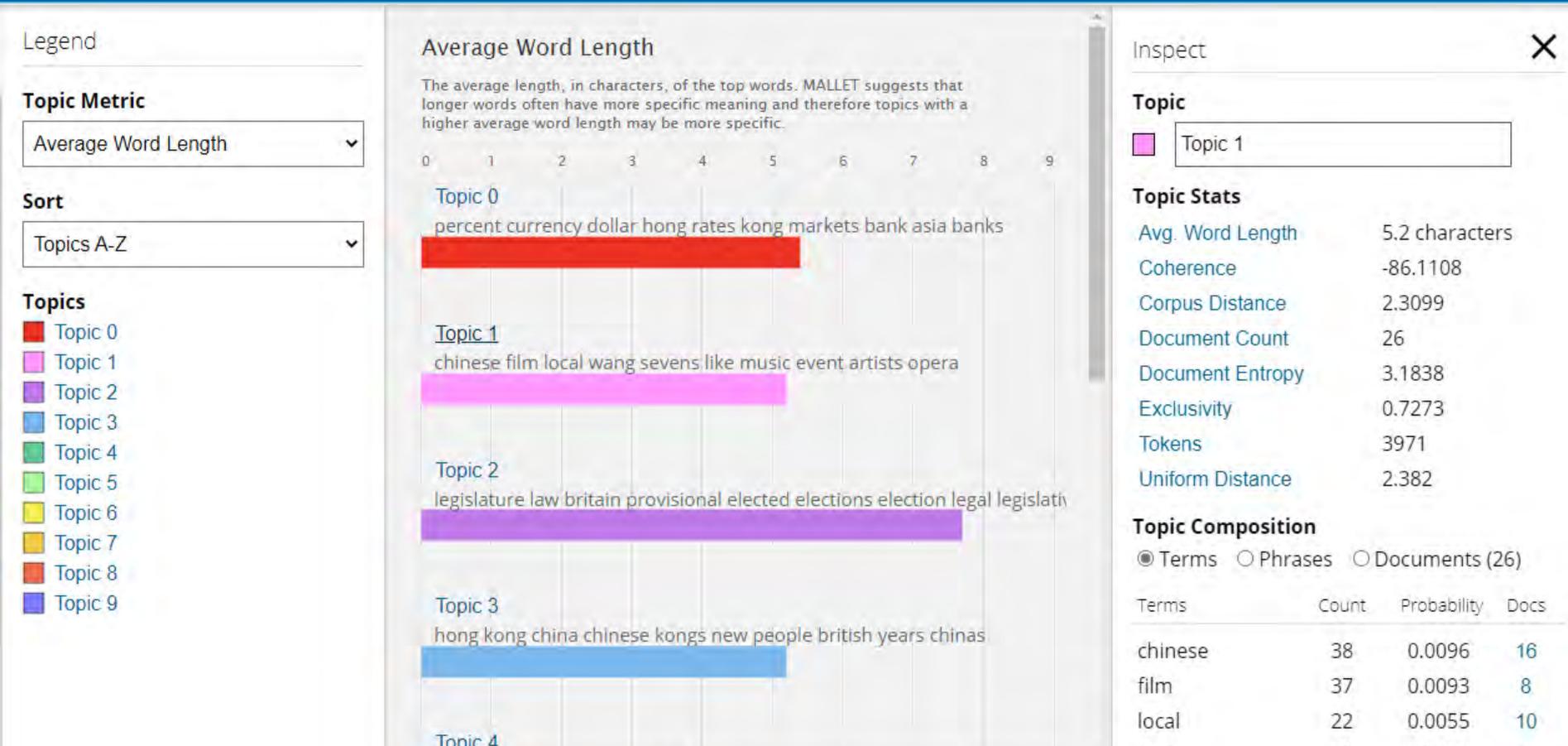
Tokens Scored : 59
Total Tokens : 495

5. Topic Modelling

- Typically used to analyze a large number of documents to find out what “**topics**” occur in common
- “Topic” = a list of keywords that occur in statistically meaningful ways. Topics are not named
- The relationship between documents and topics is a many-to-many relationship
 - Most documents contain multiple topics, and most topics appear in multiple documents, but with different weights assigned



5. Topic Modeling Output (Topic Proportion by Document)



6. Parts-of-Speech Tagger

The Parts-of-Speech tags **grammatical properties** (ex. noun, verb, adjective, adverb, etc.) to the words in documents.

It provides users with the building blocks for looking at how phrases are constructed within each document in a content set.

This tool effectively creates a lexicographical index or dictionary of a content set. In this implementation of Parts of Speech Tagger you may review how authors use of speech varies over time.



6. Parts-of-Speech Tagger Output

Legend

Parts of Speech

Click to toggle categories on and off

NOUN	INTJ	PART
PROPN	CCONJ	PUNCT
PRON	SCONJ	NUM
VERB	DET	SYM
ADV	ADP	SPACE
ADJ	AUX	OTHER

Group Data By

Author

2 of 2 selected

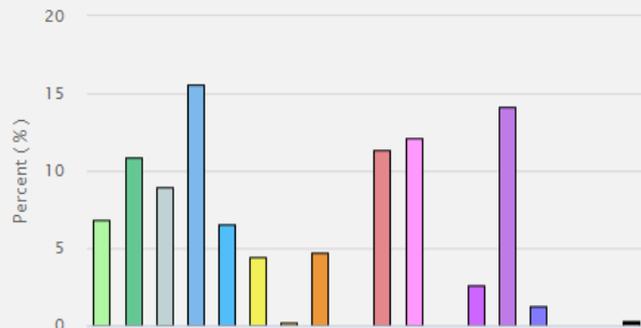
[Change selection](#)

Graphs Shown: 10

Grid Columns: 2

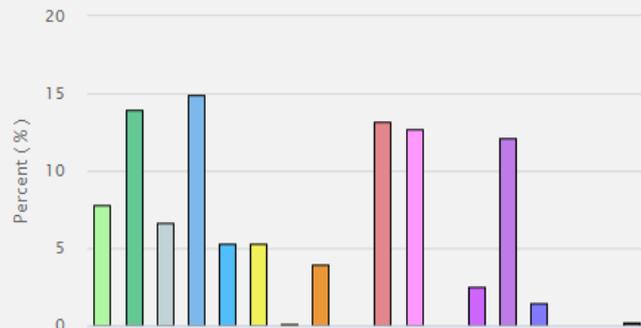
Daniel Defoe

1 documents | Tagged 145,155 Parts of Speech



Jonathan Swift

1 documents | Tagged 60,063 Parts of Speech



OCR Cleaning Configurations (Optional)

- The Lab has a built-in **OCR cleaning tool**
- Users may **change the default options** by editing stop words, removing special characters, replacing specific words, etc.
- Users can **test-run** their cleaning configurations
- The configurations can be **saved** in the user's account



Stop Words

- Stop words are words that are filtered out because they are too common and do not add much information
- Examples in English are “the” “a” “an” “at” “in” etc., etc.
- **Stop words can be edited by the user** to remove/include more words to meet their purposes
- We currently provide standard Stop Words in 26 different languages including English

STOP WORDS

[LEARN MORE](#)

Ignore stop words case

[Choose a Starter List](#) [Clear All](#)

a
about
above
across
after
afterwards
again
against
all
almost
alone
along
already
also
although
always
am
among
amongst
amongst
amount
an
and
another
any
anyhow

Why Clean?

- OCR text can contain considerable “noise”



Document Text

OCR Confidence: 73% [Learn how this text was created](#)

The Language War in Hong Kong Cantonese, Mandarin and English Compete for a Place at School

it m Hi ¥ I m .-A m A g ; * m ¥ I :li# ■.-Ni f itf % ^vf Wv yvX^AxxA : "¥ f f : x A I I m I si I* ill 3* 7 si WXvSos SSSKS.'^ ^ *-a? IN i M M i \ ' 1 ■N m M: •41' •Jg m J m > . M |\$f\$§|§ 'XV III¥. k'&'S&Xs* sSS\\|>x: \si X ill 'xxx ' ifisk ¥ ' *1 X :I KM

A student calling out his answer during a class at the Shatin Tsung Tsin Secondary School in Hong Kong. 'A, ' Jenni Meili Lau/The Washington Post
By Steven Knipp HONG KONG — For years this former British colony has had the dubious distinction of being a place where it was a serious social stigma not to be able to speak a foreign language fluently. For the city's 6 million Cantonese-speaking majority, that language was English. British rule ended in July, but the people of Hong Kong face a new quandary. Which language should their children learn? Cantonese — the local dialect spoken at home and in the street? Mandarin — China's official language? Or English — the language of international business? Whatever path the government takes, it will have far-reaching effects on a generation of young people. The language issue here is wrapped in conflicting influences and emotions, not the least of which is the local Cantonese pride that is resistant to Beijing's hopes of integrating Hong Kong into the "motherland" as quickly as possible, residents, some 60 million people in Guangdong Province nearby speak Cantonese. But Beijing considers the southern tongue just another marginal dialect among China's 8 billion people. Mandarin-speaking

Demonstration

Gale Digital Scholar Lab: The Workflow

 Build



 Clean



 Analyze

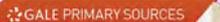
- Search
- Select and Curate
- Create Content Set (Max 10,000 items per content set)
- View Complete Search History
- Duplicate/Merge Content Sets
- Alternatively, upload your own text

- Apply default cleaning options
- Or, set your own cleaning configuration
- Edit stop words, remove special characters, etc.
- Save your cleaning configurations for later use

- Run analyses on the fly
- Six Tools for different purposes
- Download visual and tabular output
- Download OCR for max 5,000 items per content set
- Run history

Which Gale Primary Sources are available to HKBU?

Gale Primary Sources Available to HKBU

Newspapers	 THE TIMES DIGITAL ARCHIVE	 THE SUNDAY TIMES HISTORICAL ARCHIVE	 INTERNATIONAL HERALD TRIBUNE	 BRITISH LIBRARY NEWSPAPERS	 THE ILLUSTRATED LONDON NEWS	
Periodicals	 THE ECONOMIST HISTORICAL ARCHIVE	 PUNCH HISTORICAL ARCHIVE	 PICTURE POST HISTORICAL ARCHIVE	 THE TIMES LITERARY SUPPLEMENT HISTORICAL ARCHIVE	 CHINA AND THE MODERN WORLD Pt. 1	 NATIONAL GEOGRAPHIC VIRTUAL LIBRARY
Political History	 NINETEENTH CENTURY COLLECTIONS ONLINE	 ARCHIVES UNBOUND	 U.S. DECLASSIFIED DOCUMENTS ONLINE	 CHINA AND THE MODERN WORLD Pt 4, 5		
Gender History	 ARCHIVES OF SEXUALITY AND GENDER	 WOMEN'S STUDIES ARCHIVE		HTR text available	Not available in DS Lab	

Newspapers



Top UK quality paper with longest history and enormous influence, Years 1785-2014



Originally unrelated to *The Times*; analysis, investigative reporting, and culture, 1822-2006



US-owned global paper based in Paris with international readership, 1887-2013



The greatest illustrated weekly newspaper in the UK, 1842-2003



Collections of 70+ English newspapers from the 19th century, 1801-1900

Periodicals



Top UK business weekly with global readership among business elites, Years 1843-2015



The greatest satirical weekly magazine in the UK, 1841-1992



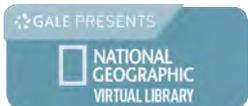
Short-lived but important UK weekly photojournalistic magazine, 1938-1957



Aka "TLS" - famous and influential weekly literary review, 1902-2013



Part 1: 17 English-language missionary, sinology, literary periodicals from China, 1817-1949



Visual US monthly magazine on nature, peoples and science, 1888-1994 *

* National Geographic not available in Digital Scholar Lab

Political & Gender History



Diplomatic reports from Asia to the US/UK, missionary reports, periodicals, and others



Part 4: British Colonial Office papers CO 129 on Hong Kong, 1841-1951 *
Part 5: British Foreign Office papers FO 17 on China, 1815–1881 *

* China and the Modern World, Parts 4-5 contain HTR text



Declassified US government documents, mainly from post-WWII 20th century



14 small sized collections related to China, Japan, Southeast Asia, and missionary activities



Publications and documents related to LGBTQ movements and their oppression

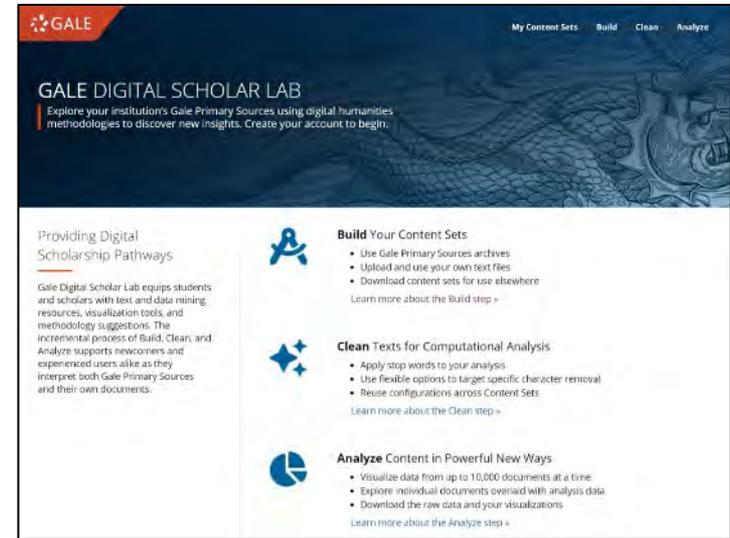


Publications and documents related to women's movements, feminism, and other topics

Gale Digital Scholar Lab

Gale's Cloud-based Text and Data Mining Research Environment

- **Access** OCR/HTR text from Gale Primary Sources
- **Build** custom content subsets
- **Clean** OCR text to reduce or control noise
- **Analyze** with simple but powerful tools
- **Export** statistical data, visualizations, OCR text in standard formats
- **Organize** and manage research



The screenshot displays the Gale Digital Scholar Lab interface. At the top, there is a navigation bar with the GALE logo and links for "My Content Sets", "Build", "Clean", and "Analyze". Below the navigation bar, the main heading reads "GALE DIGITAL SCHOLAR LAB" with a sub-heading: "Explore your institution's Gale Primary Sources using digital humanities methodologies to discover new insights. Create your account to begin." The interface is divided into three main sections, each with an icon and a brief description:

- Providing Digital Scholarship Pathways**: Accompanied by a book icon, this section describes how the lab equips students and scholars with text and data mining resources, visualization tools, and methodology suggestions. It notes that the incremental process of Build, Clean, and Analyze supports newcomers and experienced users alike as they interpret both Gale Primary Sources and their own documents.
- Build Your Content Sets**: Accompanied by a gear icon, this section lists three steps: "Use Gale Primary Sources archives", "Upload and use your own text files", and "Download content sets for use elsewhere". A link to "Learn more about the Build step" is provided.
- Clean Texts for Computational Analysis**: Accompanied by a star icon, this section lists three steps: "Apply stop words to your analysis", "Use flexible options to target specific character removal", and "Reuse configurations across Content Sets". A link to "Learn more about the Clean step" is provided.
- Analyze Content in Powerful New Ways**: Accompanied by a pie chart icon, this section lists three steps: "Visualize data from up to 10,000 documents at a time", "Explore individual documents overlaid with analysis data", and "Download the raw data and your visualizations". A link to "Learn more about the Analyze step" is provided.

Learning Center

LEARN ABOUT THE LAB

What Texts are Available?

See the full list Gale Primary Sources made available to you through your institution and learn more about each.

[View the Archives](#)

What Analyses Can I Run?

Get an overview of analysis methods, available tools, text and data mining concepts, and more.

[Explore the Analysis Tools](#)

Try Out Sample Projects

Sample Projects are a great way to acquaint yourself with the tasks you can undertake and results you can get with the Digital Scholar Lab.

[View Sample Projects](#)



[What is the Gale Digital Scholar Lab?](#)

01:16

[Transforming Scholarly Research](#)

01:10

[Supporting Faculty Research](#)

02:43

A Few Key Points that Might Help

- **Search/Retrieve before Gather/Analyze**
- **Using the Lab is a Cyclical Trial-and-Error Process**
- **Noise in OCR is to be expected**
- **The more contextual knowledge you have the better**



Thank you!