

The University of Texas at Dallas

CS 6322

Information Retrieval

Spring 2020

Class Project Report

Project TITLE: Search Engine for Countries

Group: No. 3

Students: Darrel Donald, darrel.donald@utdallas.edu

Ruchi Singh, rxs180057@utdallas.edu

Sasipreetam Morsa, sxm171330@utdallas.edu

Anshul Pardhi, anshul.pardhi@utdallas.edu

Minal Bonde, msb190000@utdallas.edu

1. Introduction

The search engine we created focuses on gathering information about Countries. The work for the different aspects of the search engine was divided as follows:

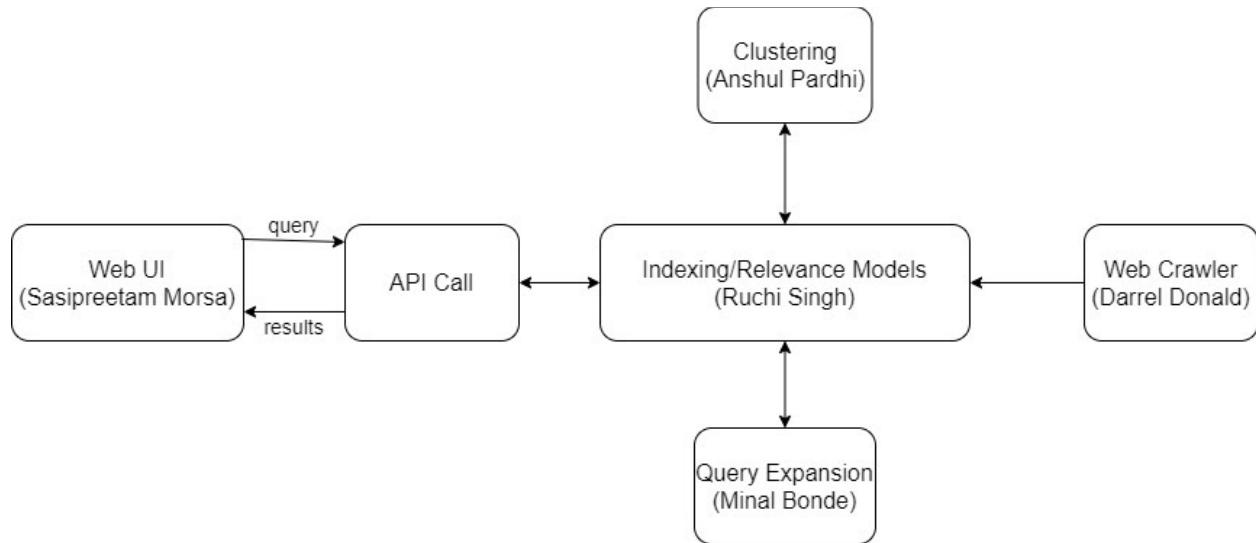
Crawling: Darrel Donald

Indexing and Relevance: Ruchi Singh

User Interface and Comparisons with Google and Bing: Sasipreetam Morsa

Clustering: Anshul Pardhi

Query Expansion and Relevance Feedback: Minal Bonde



The main thing we learned and also our biggest challenge was working with large amounts of data. When compared to the number of webpages on the internet, the amount of data we worked with is relatively small, and yet even working with this amount of data was very difficult and time consuming. As such, we also learned that there are many open-source tools available that are designed to be used with large quantities of data and that make working with this data much easier and smoother. We also created REST endpoints for transferring data between the client and the server, which was a new learning experience for some of the teammates.

Overall, we enjoyed this project because it was different from what we have seen before and gave us a new and unique experience that will help us become better developers.

2. Crawling

Initially, Scrapy was being used for crawling. Scrapy was able to obtain 100,000 unique pages, however, most of the pages were twitter links that were not useful to our search engine. We decided to try crawling again with Apache Nutch. Using Apache Nutch yielded more useful webpages.

We gave Nutch the following seed list of 50 URLs:

```
'https://en.wikipedia.org/wiki/Lists_of_countries_and_territories',
'https://en.wikipedia.org/wiki/Lists_by_country',
'https://en.wikipedia.org/wiki/List_of_sovereign_states',
'https://en.wikipedia.org/wiki/Dependent_territory',
'https://en.wikipedia.org/wiki/Country',
'https://en.wikipedia.org/wiki/Autonomous_republic',
'https://history.state.gov/countries/all',
'https://data.worldbank.org/country',
'https://www.factmonster.com/world/countries',
'https://www.infoplease.com/world/countries',
'https://www.cia.gov/library/publications/the-world-factbook/',
'https://www.traveldocs.com/world-atlas',
'https://www.countryreports.org/',
'https://www.worldometers.info/geography/how-many-countries-are-there-in-the-world/',
'https://history.state.gov/countries/all',
'https://www.countries-ofthe-world.com/all-countries.html',
'https://www.nationsonline.org/oneworld/countries_of_the_world.htm',
'https://www.britannica.com/topic/list-of-countries-1993160',
'https://worldview.stratfor.com/article/how-many-countries-are-there-world-2019',
'https://www.usnews.com/news/best-countries/slideshows/the-25-best-countries-in-the-world',
'https://www.worldatlas.com/articles/how-many-countries-are-in-the-world.html',
'https://worldpopulationreview.com/countries/how-many-countries-are-there/',
'https://theodora.com/wfb/mobile.html',
'https://www.listchallenges.com/the-195-countries-of-the-world',
```

Darrel Donald
DLD180001

'<https://www.cntraveler.com/gallery/are-these-the-best-countries-in-the-world>',
'<https://www.jetpunk.com/quizzes/how-many-countries-can-you-name>',
'<https://www.usatoday.com/story/money/2019/07/07/richest-countries-in-the-world/39630693/>',
'<https://examples.yourdictionary.com/list-of-all-countries-in-the-world.html>',
'<https://www.thoughtco.com/number-of-countries-in-the-world-1433445>',
'<https://www.boldtuesday.com/pages/alphabetical-list-of-all-countries-and-capitals-shown-on-list-of-countries-poster>',
'<https://geology.com/world/world-map.shtml>',
'<https://www.forbes.com/sites/laurabegleybloom/2020/03/20/ranked-20-happiest-countries-2020/#66d45f9c7850>',
'<https://www.businessinsider.com/best-countries-in-the-world-ranked-us-news-world-report>',
'<https://www.cdc.gov/coronavirus/2019-ncov/cases-updates/world-map.html>',
'<https://www.ethnologue.com/browse/countries>',
'<https://ourworldindata.org/countries>',
'<https://www.history.com/news/what-is-the-smallest-country-in-the-world>',
'<https://freedomhouse.org/countries/freedom-world/scores>',
'<https://www.heritage.org/index/ranking>',
'https://www.who.int/choice/demography/by_country/en/',
'<https://www.roughguides.com/gallery/most-beautiful-country-in-the-world/>',
'<https://www.weforum.org/agenda/2020/02/countries-manufacturing-trade-exports-economics/>',
'<https://www.sporcle.com/games/g/world>',
'<https://www.internetworldstats.com/list1.htm>',
'<https://www.internetworldstats.com/stats8.htm>',
'<https://www.iqair.com/us/world-most-polluted-countries>',
'<https://www.visualcapitalist.com/the-most-miserable-countries-in-the-world/>',
'<https://www.playgeography.com/games/countries-of-the-world/>',
'<https://www.ducksters.com/geography/>',
'<https://www.mapsofworld.com/map-of-countries.html>'

Darrel Donald
DLD180001

We gathered 126,434 pages useful pages using Nutch. We were aiming for more but had to cut it short due to running out of time. We ensured we did not have duplication in our crawl by using Nutch's built in remove duplicates feature. To pass on the hyperlink information, I gave read permissions to all of the data I was able to collect and explained the output of Nutch to my teammates so that they would be able to find what they needed from me and copy all of the data I collected.

3. Indexing and Relevance

Indexing

Apache Solr is used for indexing the crawled pages. The output of the crawler contains three folders, crawldb, linkdb and segments. While creating the index in Solr, all these databases have been fed as the input. The following command has been used to create the index of crawled pages in Solr with the help of nutch:

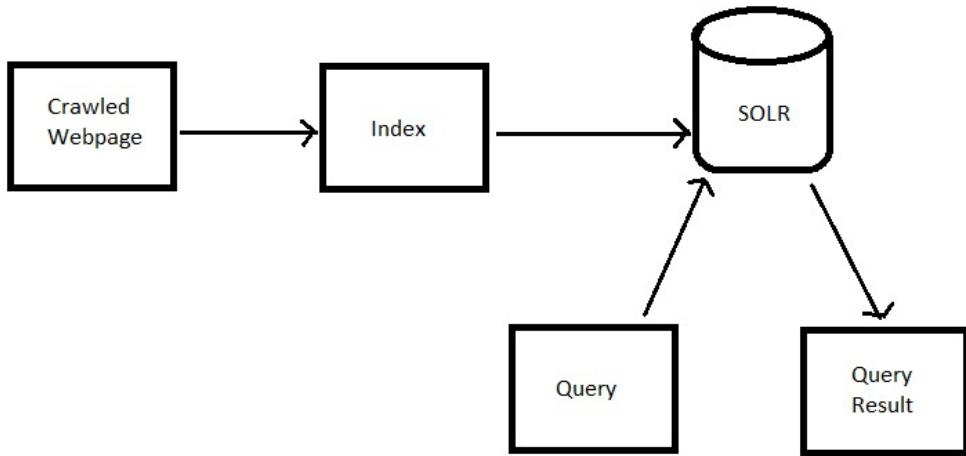
```
bin/nutch index crawl/crawldb/ -linkdb crawl/linkdb -dir crawl/segments/ -filter -normalize - deleteGone
```

bin/nutch is the binary of Apache nutch. The index command takes the content from one or more segments and passes it to all enabled IndexWriter plugins which sends the documents to Solr.

Crawl is the directory where nutch puts all the crawled pages and creates three directories mentioned above in the command.

- 1) **CrawlDb** - It contains all the links parsed by the Nutch.
- 2) **LinkDB** - It contains for each URL the outgoing and the incoming URLs.
- 3) **Segment** - It contains the list of URLs to be crawled or being crawled.

Index Creation Flow Diagram



Webgraph and Statistics

A web graph has been created from the set of crawled URLs. We have used nutch command to create the webgraph of all the crawled webpages.

```
bin/nutch webgraph -segmentDir crawl/segments/ -webgraphdb crawl/webgraphdb
```

The webgraph command of nutch takes multiple segments to process and requires an output directory in which to place the completed web graph components. The webgraph creates three different components: an inlink database, an outlink database and a node database. The inlink database is a listing of url and all of its

outlinks. The node database is a listing of url with node meta information including the number of inlinks and outlinks.

Statistics

Total Number of links = 711101

Total Number of nodes = 178112

The largest number of ingoing links = 21437

The largest number of outgoing links = 91

Web graph information connected to index

While creating index in Solr, we have run the LinkRank program of nutch to perform an iterative link analysis with the webgraph. LinkRank is a pagerank-like link analysis program that converges to stable global scores for each url. The LinkRank program starts with a common score for all urls. It then creates a global score for each url based on the number of incoming links and the scores for those links and the number of outgoing links from the page.

Relevance Models

Vector Space Relevance model

Solr provides tf-idf based relevance model to score the webpages and we made use of this model as vector space relevance model. In vector space relevance model, documents and queries are both vectors. Similarity of a document vector to a query vector can be computed by taking the cosine of the angle between them. Cosine similarity measure has been used to get the similarity between document and query. This scoring model depends on number of factors –

Term frequency - The frequency with which a term appears in a document. Given a search query, the higher the term frequency, the higher the document score.

Inverse document frequency - The rarer a term is across all documents in the index, the higher its contribution to the score.

Tf-idf weighting scheme has been used to create the relevance model. In this weighting scheme,

$$W_{d,t} = \text{tf}_{d,t} \times \text{idf}_t$$

Tf-idf weighting scheme is the most common weighting scheme in vector space relevance model.

Page Rank

To develop the relevance model based on Page Ranking, we used inbuilt nutch command. The damping factor used is 0.85. Following are the steps we followed in order to apply PageRank scores into crawled database so that query results will display web pages in the descending order of page ranking values –

1. **Create a web graph from the set of URLs fetched by the crawler**
bin/nutch webgraph -segmentDir crawl/segments/ -webgraphdb crawl/webgraphdb
2. **Run PageRank algorithm iteratively until the score converges**

```
bin/nutch linkrank -webgraphdb crawl/webgraphdb/
```

3. Update the page rank score into the crawled database

```
bin/nutch scoreupdater -crawldb crawl/crawldb - webgraphdb crawl/webgraphdb/
```

4. Finally, index the updated crawled database using Solr and nutch command:

```
bin/nutch index crawl/crawldb -linkdb crawl/linkdb -dir crawl/segments/
```

HITS (Hyperlink Induced Topic Search)

For the HITS score, we created our own python program to generate the authority and the hub score in python. HITS uses hubs and authorities to define a recursive relationship between webpages.

Given a query to a Search Engine, the set of highly relevant web pages are called **Roots**. They are potential **Authorities**. Pages which are not very relevant but point to pages in the Root are called **Hubs**. Thus, an Authority is a page that many hubs link to whereas a Hub is a page that links to many authorities.

We used the **networkx library** to create the graph through the inlinks and outlinks information of all the nodes (urls) and ran the hits algorithm on the graph to generate the hub and the authority scores. To get these inlinks, we used nutch and generated a dump of the all the inlinks of crawled pages. With the inlinks dictionary, we created the outlinks as well. Based on the query provided by the user, root set of web pages are collected based on page rank results. The root set is converted to a base set by adding out link and inlinks of the root set. After that, HIT score calculation is done by our python program and based on maximum authority score, relevant web pages are sorted and returned to the User Interface for display.

Among the webpages we crawled,

1. Highest hub score was assigned to : <https://advisor.visualcapitalist.com/> score - 0.0229952
2. Highest authority score was assigned to: <https://advisor.visualcapitalist.com/> score - 0.0229952

Topic based page ranking

For the purpose of testing Topic based page ranking, we have chosen 2 topics among countries:

1. Countries by continents
2. Countries by population

The results for **countries by continents** are as follows from our search engine: (included corresponding title and PageRank score for the webpage)

```
"url":"https://www.worldatlas.com/amp/continents.html"
"title":"7 Continents of the World - WorldAtlas.com"
"score":7.1495075 - Highest Page rank score
```

```
"title":"World Map With Continents, Map of Continents",
"url":https://www.mapsofworld.com/continents/
"score":7.141836
```

```
"title":"Continents by Population - WorldAtlas.com",
"url":https://www.worldatlas.com/articles/continents-by-population.html"
"score":7.1285152
```

"**title**":"Is North America And South America One Continent? - WorldAtlas.com",
"**url**":<https://www.worldatlas.com/articles/is-north-america-and-south-america-one-continent.html>
"**score**":7.1129417

The results for **countries by population** are as follows from our search engine: (included corresponding title and PageRank score for the webpage)

"**title**":"Population (disambiguation) - Wikipedia",
"**url**":[https://en.wikipedia.org/wiki/Population_\(disambiguation\)](https://en.wikipedia.org/wiki/Population_(disambiguation))
"**score**":4.105301 - Highest Page rank score

"**title**":"Countries With the Lowest Population Density - WorldAtlas.com",
"**url**":<https://www.worldatlas.com/amp/articles/countries-with-the-lowest-population-density.html>
"**score**":4.096262

"**title**":"List of European countries by population - Wikipedia",
"**url**":https://en.wikipedia.org/wiki/List_of_European_countries_by_population
"**score**":4.090098

"**title**":"List of Asian countries by population - Wikipedia",
"**url**":https://en.wikipedia.org/wiki/List_of_Asian_countries_by_population
"**score**":4.085929

Collaboration with UI and test relevance models results

I have generated several queries to test the various relevance models I have implemented in this search engine. I have generated 5 queries and tested on the UI by differentiating the results of pagerank and HITS relevance model results. I have judged the results in some ways such as, HITS relevance model was giving most relevant results in comparison of page rank. For example,

For query – india

HITS gives the results in which the pages should come first as Wikipedia but page rank ranked the other pages more than Wikipedia. In the 1st screenshot, HITS option is selected.

The screenshot shows a web browser window with three tabs, all titled "Countries Search Engine". The current tab displays search results for the query "india".

Search Query: india

Relevance Model Options:

- Page Rank
- HITS

Clustering Options:

- Flat Clustering
- Hierarchical Clustering

Query Expansion Option:

- Association
- Metric
- Scalar

Search

Results:

- [India - Wikipedia](https://en.wikipedia.org/wiki/India)
https://en.wikipedia.org/wiki/India
India Wikipedia India From Wikipedia the free encyclopedia Jump to navigation Jump to search This article is about the Republic of India For other uses see India disambiguation Bharat redir
- [India Map, Map of India](#)
https://www.mapsofworld.com/india/
India Map Map of India Menu World Map World Maps Political Map of the World Physical Map of the World Blank World Map World Map for Kids Earth Map World Atlas World Time Zones Map Counties in UK
- [Sino-Tibetan | Ethnologue](#)
https://www.ethnologue.com/subgroups/sino-tibetan
Sino Tibetan Ethnologue Skip to main content Login Shopping cart Ethnologue Languages Countries Guides About Plans Pricing Sino Tibetan Print Sino Tibetan Chinese Chinese Gan gan A
- [Indo-European | Ethnologue](#)
https://www.ethnologue.com/subgroups/indo-european
Indo European Ethnologue Skip to main content Login Shopping cart Ethnologue Languages Countries Guides About Plans Pricing Indo European Print Indo European Albanian Gheg Albanian
- [Political Map of India](#)

In the 2nd screenshot, page rank option is selected.

india

Relevance Model Options:

Page Rank HITS

Clustering Options:

Flat Clustering Hierarchical Clustering

Query Expansion Option:

Association Metric Scalar

Search

[India Map, Map of India](#)
https://www.mapsofworld.com/india/
India Map Map of India Menu World Map World Maps Political Map of the World Physical Map of the World Blank World Map World Map for Kids Earth Map World Atlas World Time Zones Map Counties in UK

[Sino-Tibetan | Ethnologue](#)
https://www.ethnologue.com/subgroups/sino-tibetan
Sino Tibetan Ethnologue Skip to main content Login Shopping cart Ethnologue Languages Countries Guides About Plans Pricing Sino Tibetan Print Sino Tibetan Chinese Chinese Gan gan A

[Indo-European | Ethnologue](#)
https://www.ethnologue.com/subgroups/indo-european
Indo European Ethnologue Skip to main content Login Shopping cart Ethnologue Languages Countries Guides About Plans Pricing Indo European Print Indo European Albanian Gheg Albanian

[Political Map of India](#)
https://www.mapsofworld.com/india/india-political-map.html
Political Map of India Menu World Map World Maps Political Map of the World Physical Map of the World Blank World Map World Map for Kids Earth Map World Atlas World Time Zones Map Counties in UK

[India \(disambiguation\) - Wikipedia](#)

For query – hong kongIn the 1st screenshot, HITS option is selected. Wikipedia page is coming on top.

Countries Search Engine

hong kong

Relevance Model Options:

Page Rank HITS

Clustering Options:

Flat Clustering Hierarchical Clustering

Query Expansion Option:

Association Metric Scalar

Search

Hong Kong - Wikipedia
https://en.wikipedia.org/wiki/Hong_Kong
 Hong Kong Wikipedia Hong Kong From Wikipedia the free encyclopedia Jump to navigation Jump to search Special administrative region of China HK redirects here For other uses see Hong Kong disambiguation

History of Hong Kong - Wikipedia
https://en.wikipedia.org/wiki/History_of_Hong_Kong
 History of Hong Kong Wikipedia History of Hong Kong From Wikipedia the free encyclopedia Jump to navigation Jump to search This article is part of a series on the History of Hong Kong Timeline Prehistory

South Asians in Hong Kong - Wikipedia
https://en.wikipedia.org/wiki/South_Asians_in_Hong_Kong
 South Asians in Hong Kong Wikipedia South Asians in Hong Kong From Wikipedia the free encyclopedia Jump to navigation Jump to search South Asians in Hong Kong Total population More than

Hong Kong City Map
<https://www.mapsofworld.com/hong-kong/hong-kong-city-map.html>
 Hong Kong City Map Menu World Map World Maps Political Map of the World Physical Map of the World Blank World Map World Map for Kids Earth Map World Atlas World Time Zones Map Counties in UK Different

In the 2nd screenshot, page rank option is selected. History of hong kong page is coming on top.

Countries Search Engine

hong kong

Relevance Model Options:

Page Rank HITS

Clustering Options:

Flat Clustering Hierarchical Clustering

Query Expansion Option:

Association Metric Scalar

Search

History of Hong Kong - Wikipedia
https://en.wikipedia.org/wiki/History_of_Hong_Kong
 History of Hong Kong Wikipedia History of Hong Kong From Wikipedia the free encyclopedia Jump to navigation Jump to search This article is part of a series on the History of Hong Kong Timeline Prehistory

South Asians in Hong Kong - Wikipedia
https://en.wikipedia.org/wiki/South_Asians_in_Hong_Kong
 South Asians in Hong Kong Wikipedia South Asians in Hong Kong From Wikipedia the free encyclopedia Jump to navigation Jump to search South Asians in Hong Kong Total population More than

Hong Kong City Map
<https://www.mapsofworld.com/hong-kong/hong-kong-city-map.html>
 Hong Kong City Map Menu World Map World Maps Political Map of the World Physical Map of the World Blank World Map World Map for Kids Earth Map World Atlas World Time Zones Map Counties in UK Different

COVID-19 Information | U.S. Consulate General Hong Kong & Macau
<https://hk.usconsulate.gov/covid-19-information/>
 COVID-19 Information U.S. Consulate General Hong Kong Macau English U.S. Consulate General Hong Kong Macau Social Search Facebook Instagram Twitter Search Global Level Health Advisory

For query – Malaysia population

In the 1st screenshot, page rank option is selected. Demographics page is coming on top.

The screenshot shows a web browser window titled "Countries Search Engine" with three tabs open, all showing the same search results for "malaysia population". The search bar at the top contains "malaysia population". Below the search bar are sections for "Relevance Model Options" (with "Page Rank" selected), "Clustering Options" (with "Flat Clustering" selected), and "Query Expansion Option" (with "Association" selected). A "Search" button is present. The main content area displays four search results:

- Malaysia Population 2020 (Demographics, Maps, Graphs)**
https://worldpopulationreview.com/countries/malaysia-population/
Malaysia Population Demographics Maps Graphs World Population Review Mobile Navigation Home Continents Countries World Cities US States US Counties US Cities Zips Canadian Provinces Coronavir
- Malaysia Demographics 2020 (Population, Age, Sex, Trends) - Worldometer**
https://www.worldometers.info/demographics/malaysia-demographics/
Malaysia Demographics Population Age Sex Trends Worldometer Coronavirus Population W Demographics Malaysia Demographics Malaysia Demographics Population Median Age Dependency Ratio Fertil
- Malaysia Population (2020) - Worldometer**
https://www.worldometers.info/world-population/malaysia-population/
Malaysia Population Worldometer Coronavirus Population W Population World Asia South Eastern Asia Malaysia Malaysia Population LIVE retrieving data Malaysia Population Year
- Demographics of Malaysia - Wikipedia**
https://en.wikipedia.org/wiki/Demographics_of_Malaysia
Demographics of Malaysia Wikipedia Demographics of Malaysia From Wikipedia the free encyclopedia Jump to navigation Jump to search Demographics of Malaysia Indicator Rank Measure Economy GDP PPP

In the 2st screenshot, HITS option is selected. Worldometer page is coming on top about population.

Countries Search Engine

malaysia population

Relevance Model Options:

- Page Rank HITS

Clustering Options:

- Flat Clustering Hierarchical Clustering

Query Expansion Option:

- Association Metric Scalar

[Search](#)

[Malaysia Population \(2020\) - Worldometer](#)
<https://www.worldometers.info/world-population/malaysia-population/>
 Malaysia Population Worldometer Coronavirus Population W Population World Asia South Eastern Asia Malaysia Malaysia Population LIVE retrieving data Malaysia Population Year

[Malaysia Population 2020 \(Demographics, Maps, Graphs\)](#)
<https://worldpopulationreview.com/countries/malaysia-population/>
 Malaysia Population Demographics Maps Graphs World Population Review Mobile Navigation Home Continents Countries World Cities US States US Counties US Cities Zips Canadian Provinces Coronavir

[Malaysia Demographics 2020 \(Population, Age, Sex, Trends\) - Worldometer](#)
<https://www.worldometers.info/demographics/malaysia-demographics/>
 Malaysia Demographics Population Age Sex Trends Worldometer Coronavirus Population W Demographics Malaysia Demographics Malaysia Demographics Population Median Age Dependency Ratio Fertil

[Demographics of Malaysia - Wikipedia](#)
https://en.wikipedia.org/wiki/Demographics_of_Malaysia
 Demographics of Malaysia Wikipedia Demographics of Malaysia From Wikipedia the free encyclopedia Jump to navigation Jump to search Demographics of Malaysia Indicator Rank Measure Economy GDP PPP

For Query – Philippines

Page rank is selected and it gives geography pages on top which are more relevant.

Countries Search Engine

Philippines

Relevance Model Options:

- Page Rank HITS

Clustering Options:

- Flat Clustering Hierarchical Clustering

Query Expansion Option:

- Association Metric Scalar

[Search](#)

[Geography and Fact Sheet About the Philippines](#)
<https://www.thoughtco.com/geography-of-the-philippines-1435646>
 Geography and Fact Sheet About the Philippines Menu Home The Philippines Geography and Fact Sheet Search Search the site GO Geography Country Information Basics Physical Geography Political Geography

[Where is Philippines? | Where is Philippines in the World?](#)
<https://worldpopulationreview.com/country-locations/where-is-philippines/>
 Where is Philippines Where is Philippines in the World World Population Review Mobile Navigation Home Continents Countries World Cities US States US Counties US Cities Zips Canadian Provinces Coro

[Map of the Philippines](#)
<https://www.mapsofworld.com/philippines/>
 Map of the Philippines Menu World Map World Maps Political Map of the World Physical Map of the World Blank World Map World Map for Kids Earth Map World Atlas World Time Zones Map Counties in UK D

[Philippines Population 2020 \(Demographics, Maps, Graphs\)](#)
<https://worldpopulationreview.com/countries/philippines-population/>
 Philippines Population Demographics Maps Graphs World Population Review Mobile Navigation Home Continents Countries World Cities US States US Counties US Cities Zips Canadian Provinces Corona

HITS is selected and it gives less relevant results for this query in comparison of page rank.

The screenshot shows a web browser window with four tabs open. The active tab is titled "Philippines - Google Search" and displays search results for the query "Philippines". The results include links to "Philippines Population (2020) - Worldometer", "Philippines Map", "Geography and Fact Sheet About the Philippines", and "Where is Philippines? | Where is Philippines in the World?". The browser interface includes a navigation bar with back, forward, and search buttons, and a status bar at the bottom.

Relevance Model Options:
 Page Rank HITS

Clustering Options:
 Flat Clustering Hierarchical Clustering

Query Expansion Option:
 Association Metric Scalar

Search

Philippines Population (2020) - Worldometer
https://www.worldometers.info/world-population/philippines-population/
Philippines Population Worldometer Coronavirus Population W Population World Asia South Eastern Asia Philippines Philippines Population LIVE retrieving data Philippines Population

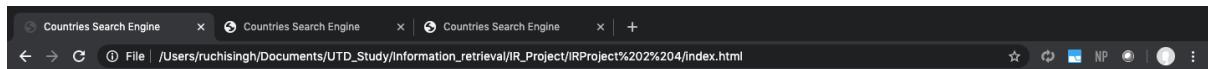
Philippines Map
https://www.worldatlas.com/maps/philippines.html
Philippines Map Trending Covid Cases And Deaths Per State In The US Could There Be Another Great Depression Celebrities Who Have Recovered From COVID World Map Asia Philippines Maps Of P

Geography and Fact Sheet About the Philippines
https://www.thoughtco.com/geography-of-the-philippines-1435646
Geography and Fact Sheet About the Philippines Menu Home The Philippines Geography and Fact Sheet Search Search the site GO Geography Country Information Basics Physical Geography Political Geography

Where is Philippines? | Where is Philippines in the World?
https://worldpopulationreview.com/country-locations/where-is-philippines/
Where is Philippines Where is Philippines in the World World Population Review Mobile Navigation Home Continents Countries World Cities US States US Counties US Cities Zips Canadian Provinces Coro

For Query – Zimbabwe

Page rank is selected and it gives the same results as HITS.



Countries Search Engine

zimbabwe

Relevance Model Options:

Page Rank HITS

Clustering Options:

Flat Clustering Hierarchical Clustering

Query Expansion Option:

Association Metric Scalar

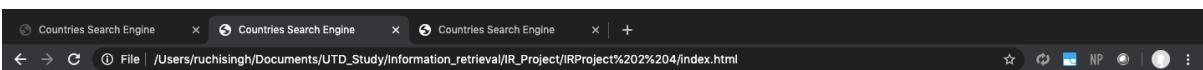
Search

[Zimbabwe -- Home - allAfrica.com](#)
<https://allafrica.com/zimbabwe/>
 Zimbabwe Home allAfrica.com Search Countries All Countries Algeria Angola Benin Botswana Burkina Faso Burundi Cameroon Cape Verde Central African Republic Chad Comoros Congo Brazzaville Congo Kin

[Zimbabwe Population 2020 \(Demographics, Maps, Graphs\)](#)
<https://worldpopulationreview.com/countries/zimbabwe-population/>
 Zimbabwe Population Demographics Maps Graphs World Population Review Mobile Navigation Home Continents Countries World Cities US States US Counties US Cities Zips Canadian Provinces Coronavir

[Zimbabwe - Wikipedia](#)
<https://en.wikipedia.org/wiki/Zimbabwe>
 Zimbabwe Wikipedia Zimbabwe From Wikipedia the free encyclopedia This is the latest accepted revision reviewed on April Jump to navigation Jump to search ZWE redirects here For other

[Zimbabwe Map](#)
<https://www.mapsofworld.com/zimbabwe/>
 Zimbabwe Map Maps of World Current Credible Consistent Search Maps World Maps Political Map of the World Physical Map of the World Blank World Map World Map for Kids Earth Map World Atlas World Ti



Countries Search Engine

zimbabwe

Relevance Model Options:

Page Rank HITS

Clustering Options:

Flat Clustering Hierarchical Clustering

Query Expansion Option:

Association Metric Scalar

Search

[Zimbabwe Population \(2020\) - Worldometer](#)
<https://www.worldometers.info/world-population/zimbabwe-population/>
 Zimbabwe Population Worldometer Coronavirus Population W Population World Africa Eastern Africa Zimbabwe Zimbabwe Population LIVE retrieving data Zimbabwe Population Yearly

[Zimbabwe -- Home - allAfrica.com](#)
<https://allafrica.com/zimbabwe/>
 Zimbabwe Home allAfrica.com Search Countries All Countries Algeria Angola Benin Botswana Burkina Faso Burundi Cameroon Cape Verde Central African Republic Chad Comoros Congo Brazzaville Congo Kin

[Zimbabwe Population 2020 \(Demographics, Maps, Graphs\)](#)
<https://worldpopulationreview.com/countries/zimbabwe-population/>
 Zimbabwe Population Demographics Maps Graphs World Population Review Mobile Navigation Home Continents Countries World Cities US States US Counties US Cities Zips Canadian Provinces Coronavir

[Zimbabwe - Wikipedia](#)
<https://en.wikipedia.org/wiki/Zimbabwe>
 Zimbabwe Wikipedia Zimbabwe From Wikipedia the free encyclopedia This is the latest accepted revision reviewed on April Jump to navigation Jump to search ZWE redirects here For other

[Zimbabwe Map](#)

Collaboration with Clustering

I have collaborated with clustering to improve the results of the relevance models I have used. For all above queries mentioned, clustering (Hierarchical clustering) gives better results in comparison of page rank and HITS relevance models.

User Interface and Comparisons with Google and Bing

4. User Interface and Comparisons with Google and Bing (Sasipreetam Morsa)

The visual part of the User Interface is made of a form and 3 iframes.

The form includes a search bar to enter queries, radio buttons to select the method used to rank the results, and a submit button that, when clicked, runs JavaScript functions to display the results in the different iframes.

The 3 iframes contain our custom search engine, Google, and Bing, which are displayed in that order. The Google and Bing iframes start with their respective home pages, while our custom iframe starts out being empty.

To get the results for the user interface, I make a get request, with the query and type as parameters, to an API that connected to the relevance models, provided by Ruchi. The API uses the query and type parameters to build a JSON file with the results ordered by ranking. The different type parameters are as follow:

- Page Rank – API returns JSON with page rank based relevance model
- HITS – API returns JSON with HITS based relevance model
- Flat Clustering – API returns JSON with Flat Clustering based relevance model
- Hierarchical Clustering – API returns JSON with Hierarchical Clustering based relevance model
- Metric Query Expansion – API returns JSON with relevance model based on metric cluster query expansion
- Association Query Expansion – API returns JSON with relevance model based on metric association query expansion
- Scalar Query Expansion – API returns JSON with relevance model based on metric scalar query expansion

Ruchi and I together only used a few queries to test the communication between the User Interface and Ruchi's API, but separately I have tested over 20 queries to both ensure that the API requests could handle many requests back to back and to check that each of the different parameters were being sent to the API properly.

The work for clustering was done by Anshul. For the given query, Ruchi's relevance model would provide the top 50 results to Anshul's clustering program. The clustering program would then rearrange the results such that the results would start with the same top results but would then have the results in the

User Interface and Comparisons with Google and Bing

same cluster as the top result be the next results. So basically, the results were rearrange based on clusters containing the top result and so on.

Overall, when compared to large search engines like Google or Bing, our search engine is not as reliable. While there are some queries on which our search engine can keep up with or even beat the large search engines, it is not on the same level as Google or Bing because we don't have access to as many webpages as those search engines and we also don't have a way to keep updating the results in real time based on which websites more users are gravitating towards given a certain type of query.

For the demonstration, we have chosen queries that we believe have been showing both accurate results and show improvement as we include clustering or query expansion.

Below are the results of some of the queries we have tried:

Results for query "costa rica":

Countries Search Engine

costa rica

Relevance Model Options:
 Page Rank HITS

Clustering Options:
 Flat Clustering Hierarchical Clustering

Query Expansion Option:
 Association Metric Scalar

[File:Coat of arms of Costa Rica.svg - Wikipedia](#)
https://en.wikipedia.org/wiki/File:Coat_of_arms_of_Costa_Rica.svg
 File Coat of arms of Costa Rica svg Wikipedia File Coat of arms of Costa Rica svg From Wikipedia the free encyclopedia Jump to navigation Jump to search File File history File usage Global file usa

[Costa Rica Map, Map of Costa Rica](#)
<https://www.mapsofworld.com/costa-rica/>
 Costa Rica Map Map of Costa Rica Menu World Map World Maps Political Map of the World Physical Map of the World Blank World Map World Map for Kids Earth Map World Atlas World Time Zones Map Count

[Costa Rica - Wikipedia](#)
https://en.wikipedia.org/wiki/List_of_Costa_Ricans
 Costa Rica Wikipedia Costa Rica From Wikipedia the free encyclopedia Redirected from List of Costa Ricans Jump to navigation Jump to search For cities in other countries see Costa Rica Sinal

[Costa Rica - Wikipedia](#)
https://en.wikipedia.org/wiki/Costa_Rica
 Costa Rica Wikipedia Costa Rica From Wikipedia the free encyclopedia Jump to navigation Jump to search For cities in other countries see Costa Rica Sinaloa and Costa Rica Mato Grosso do Sul Fo

User Interface and Comparisons with Google and Bing

Google Search

Google costa rica

All Maps News Images Videos More Settings Tools

About 868,000,000 results (1.05 seconds)

[www.visitcostarica.com](#) ... Visit Costa Rica | Costa Rica Tourism Official website of Costa Rica. Here, visitors enjoy lovely tropical beaches, grand adventures, the wonders of nature and scintillating culture - all the necessary components of an ... About Costa Rica - Costa Rica | Visit Costa Rica - Costa Rica Maps - Contact Us

People also ask

- Is it dangerous to visit Costa Rica?
- What is the best time of year to go to Costa Rica?
- Is Costa Rica part of USA?




Map of Costa Rica
Map data ©2020 INEGI

Costa Rica

Country in Central America

Costa Rica is a rugged, rainforested Central American country with coastlines on the Caribbean and Pacific. Though its capital, San Jose, is home to cultural

Google costa rica

en.wikipedia.org › wiki › Costa_Rica Costa Rica - Wikipedia

Costa Rica officially the Republic of Costa Rica (Spanish: República de Costa Rica), is a country in Central America, bordered by Nicaragua to the north, the ... Capital and largest city: San José; 9°56'N 84... Calling code: +506 Government: Unitary presidential constitutional ... Demonym(s): Costa Rican; Tico(a) History of Costa Rica - San José, Costa Rica - Geography of Costa Rica - Limón

Top stories



Costa Rica's economy will shrink 3.6% in 2020 due to COVID-19, according to Central Bank
The Tico Times · 1 day ago



The latest Costa Rica flight information – The Tico Times | Costa Rica News | Travel | Real Estate
The Tico Times · 2 days ago

spider monkeys and quetzal birds.

Capital: San José

Currency: Costa Rican colón

Recognized regional languages: Mekatelyu; Bri bri; Patois

Population: 4.99 million (2018) World Bank

Continent: North America

Official language: Spanish

People also search for View 15+ more

Panama
Puerto Rico
Cuba
Belize

Feedback

Google costa rica

www.cia.gov › publications › the-world-factbook › geos Central America :: Costa Rica – The World Factbook - Central ...

Central America :: Costa Rica Print. Page last updated on March 13, 2020. Costa Rica Flag. The World Factbook Country/Location Flag Modal x.

www.lonelyplanet.com › costa-rica Costa Rica travel | Central America - Lonely Planet

Explore Costa Rica holidays and discover the best time and places to visit. | Centering yourself on a surfboard or yoga mat, descending into bat-filled caves or ...

www.anywhere.com › costa-rica Costa Rica Travel Guide 2020 | Anywhere.com

Plan the ultimate vacation in Costa Rica with our 2020 Travel Guide to 2020, a travel destination perfect for families, couples and more.

★★★★★ Rating: 5 - 436 reviews

www.travelandleisure.com › Trip Ideas 12 Reasons Why Costa Rica Is One of the Best Vacation

User Interface and Comparisons with Google and Bing

Bing Search

costa rica

ALL NEWS IMAGES VIDEOS MAPS SHOPPING

36,100,000 Results Any time ▾

News about Costa Rica
bing.com/news

Costa Rica pushes back coronavirus, reducing current cases
While parts of Latin America enter the toughest phase of the coronavirus pandemic, Costa Rica has for the ...
Reuters on MSN.com · 2d

Costa Rica's Fonatet project making continues to make progress
Costa Rica's Fonatet project making continues to make progress
Ondrej Avendaño - AP

CINDE announces: East West invests US\$2 million in Costa Rica and will be hiring 50 more people
CINDE announces: East West invests US\$2 million in Costa Rica and will be hiring 50 more people

Costa Rica
Country in Central America
Costa Rica, officially the Republic of Costa Rica, is a

Bing Search

costa rica

ALL IMAGES VIDEOS MAPS NEWS SHOPPING

What would you like to know about this country?
things to do currency time earthquake flag

Visit Costa Rica | Costa Rica Tourism Official website
<https://www.visitcostarica.com> ▾
What makes Costa Rica Unique? What makes Costa Rica Unique? components of an ideal vacation or holiday; authentic experiences. Fill your calendar with the most wonderful natural and wildlife events. Costa Rica occupies a privileged spot with beaches in the Caribbean Sea and the Pacific Ocean. Costa Rica is a land of volcanoes, rainforests and ...

Things To Do
In Costa Rica will find plenty of exhilarating ...
Culture · Cruises · Families

Costa Rica
Costa Rica is one of the most biodiverse ...
Visit Costa Rica - Faq's

See more ▾

Costa Rica
Country in Central America

Bing Search

Costa Rica - Wikipedia
https://en.wikipedia.org/Wiki/Costa_Rica ▾

National Park	Antonio National P...	National Park	Volcano National P...	National Park
Carlos Alvarado Quesada President	Laura Chinchilla Barr Former President	Epsy Campbell Barr Vice President	Luis Guillermo Solís Former Pres...	Oscar Arias Former Pres...

Costa Rica , officially the Republic of Costa Rica (Spanish: República de Costa Rica), is a country in Central America, bordered by Nicaragua to the north, the Caribbean Sea to the northeast, Panama to the southeast, the Pacific Ocean to the southwest, and Ecuador to the south of Cocos Island. It has a population of around 5 million in a land area of 51,060 square kilometers (19,714 square miles). An estimated 333,980 people live in the capital and largest city, San José, with around 2 million people in the surrounding metropolitan area. +

Wikipedia - Text under CC-BY-SA license

Costa Rica 2020: Best of Costa Rica Tourism - Tripadvisor
https://www.tripadvisor.com/Tourism-g291982-Costa_Rica-Vacations.html ▾
Costa Rica Embraced by the Pacific Ocean and the Caribbean Sea, Costa Rica has been attracting eco-tourists for decades, thanks to rich biodiversity and abundant wildlife. Here you'll find splendid beaches that spill into sparkling blue waters, peaceful conservation areas, and outstanding adventure sports.

Videos of costa rica
bing.com/videos

User Interface and Comparisons with Google and Bing

Results for query "israel":

Countries Search Engine

israel

Relevance Model Options:
 Page Rank HITS

Clustering Options:
 Flat Clustering Hierarchical Clustering

Query Expansion Option:
 Association Metric Scalar

[Bene Israel - Wikipedia](#)
https://en.wikipedia.org/wiki/Bene_Israel
Bene Israel Wikipedia Bene Israel From Wikipedia the free encyclopedia Jump to navigation Jump to search Jewish ethnic group Not to be confused with Bani Isra il or Beta Israel Bene Israel

[Israel Population 2020 \(Demographics, Maps, Graphs\)](#)
https://worldpopulationreview.com/countries/israel-population/
Israel Population Demographics Maps Graphs World Population Review Mobile Navigation Home Continents Countries World Cities US States US Counties US Cities Zips Canadian Provinces Coronavirus

[Demographics of Israel - Wikipedia](#)
https://en.wikipedia.org/wiki/Demographics_of_Israel
Demographics of Israel Wikipedia Demographics of Israel From Wikipedia the free encyclopedia Jump to navigation Jump to search Demographics of Israel Population of Israel since Population

[Israel - Wikipedia](#)
https://en.wikipedia.org/wiki/Israel
Israel Wikipedia Israel From Wikipedia the free encyclopedia Jump to navigation Jump to search This article is about the modern country For other uses see Israel disambiguation Country in Wes

[Israelis - Wikipedia](#)

... more results

User Interface and Comparisons with Google and Bing

Google Search

The screenshot shows the Google search interface for the query "israel". The search bar at the top contains the word "israel". Below the search bar, there are tabs for "All", "News", "Maps", "Images", "Videos", and "More". On the right side of the search bar are "Settings" and "Tools" buttons, and a "Sign in" button. The main search results area displays "About 1,150,000,000 results (0.97 seconds)". A "Top stories" section features three news items with images and titles: "Israeli archaeologists find hidden pattern at 'world's oldest temple' Göbekli Tepe" (Haaretz, 13 hours ago), "Foreign Ministry: Egyptian TV show predicting end of Israel 'unacceptable'" (The Jerusalem Post, 10 hours ago), and "Israel to hold Memorial Day ceremonies without audience amid coronavirus restrictions" (Haaretz, 2 mins ago). To the right of these stories is a box containing the flag of Israel and a map of the Middle East region, with a caption about Israel being a Middle Eastern country on the Mediterranean Sea. The map includes labels for Jordan, Israel, and Cairo.

Google Search

This screenshot of the Google search results for "israel" shows a more detailed view. The search bar now includes the URL "en.wikipedia.org · wiki · Israel". The main content area features a summary of Israel's history, mentioning the 250,000 Palestinian Arabs who fled or were expelled in 1948, and the establishment of the State of Israel. It also provides information about its capital (Jerusalem), largest city (Jerusalem), ethnic groups (74.2% Jewish, 20.9% Arab), and currency (New shekel (ILS)). Below this summary is a "People also ask" section with questions like "What was Israel called before 1948?", "Is it safe in Israel?", and "What is special about Israel?". To the right of the main content, there is a sidebar with facts about Israel: population (8.884 million (2018)), capital (Jerusalem), date format (yyyy-mm-dd (AM); dd-mm-yyyy (CE)), continent (Asia), prime minister (Benjamin Netanyahu), and official language (Hebrew). There is also a "People also search for" section with links to "Jerusalem", "Iraq", "Syria", and "Russia", each accompanied by a small flag icon.

Google Search

This screenshot shows the Google search results for "israel" with a focus on news sources. The search bar contains "israel". The results list several news articles from different sources: 1. "Israel | Facts, History, & Map | Britannica" from www.britannica.com, describing Israel as a country in the Middle East located at the eastern end of the Mediterranean Sea. 2. "Israel - The New York Times" from www.nytimes.com, stating that Israel is seated on a portion of land in the Middle East, known from 1920 to 1948 as Palestine. 3. "The Times of Israel | News from Israel, the Middle East and the ..." from www.timesofisrael.com, described as a one-stop site for news, features, live blogs, and more. 4. "www.cia.gov · publications · the-world-factbook · geos" from www.cia.gov, likely referring to the CIA World Factbook. Each news item has a brief description and a link to the full article.

User Interface and Comparisons with Google and Bing

Bing Search

The Bing search interface for 'israel' includes a search bar with the query 'israel', a sign-in button, a user icon (200), and a menu icon. Below the search bar are filters for ALL, NEWS, IMAGES, VIDEOS, MAPS, and SHOPPING. The results show 152,000,000 results from bing.com/news. A news card for 'Israel's once-mighty Labor party weighs unity with Netanyahu' is displayed, along with other news snippets and images. To the right is a map of the Middle East with Israel highlighted, showing major cities like Tel Aviv-Yafo, Jerusalem, Amman, and Damascus.

Bing Search

The Bing search results for 'israel' include a link to the Wikipedia page for Israel. The page content discusses Israel's location, borders, and geographical features. It also includes sections for Overview, Etymology, History, Geography and environment, Demographics, and a list of former prime ministers. To the right of the main content are boxes for various Israeli landmarks and historical figures, along with a sidebar for people also asking about 'Is "Israeli" a nationality?' and a suggestion to 'Suggest an edit'.

Bing Search

The Bing search results for 'israel' also include links to the Britannica facts page, the Tripadvisor tourism page, and a video section. The Britannica page provides general information about Israel's topography and history. The Tripadvisor page highlights Israel as a destination for tourism. The video section shows three clips: 'Israel Aerospace Industries', 'OFFICIAL_Somewhere over', and 'Ethiopian Jews'. A sidebar on the right shows a preview of a mobile device displaying the search results.

User Interface and Comparisons with Google and Bing

Results for query “pakistan”:

Countries Search Engine

pakistan

Relevance Model Options:
 Page Rank HITS

Clustering Options:
 Flat Clustering Hierarchical Clustering

Query Expansion Option:
 Association Metric Scalar

[Pakistan Map, Map of Pakistan, Information and Interesting Facts of Pakistan](#)
<https://www.mapsofworld.com/pakistan/>
 Pakistan Map Map of Pakistan Information and Interesting Facts of Pakistan Menu World Map World Maps Political Map of the World Physical Map of the World Blank World Map World Map for Kids Earth

[Pakistan - WorldAtlas.com](#)
<https://www.worldatlas.com/topics/pakistan/>
 Pakistan WorldAtlas com Trending Covid Cases And Deaths Per State In The US Could There Be Another Great Depression Celebrities Who Have Recovered From COVID Autumn landscape in Chitral Gol

[Foreign relations of Pakistan - Wikipedia](#)
https://en.wikipedia.org/wiki/Foreign_relations_of_Pakistan
 Foreign relations of Pakistan Wikipedia Foreign relations of Pakistan From Wikipedia the free encyclopedia Jump to navigation Jump to search The flag of Pakistan Pakistan This article is part of a

[East Pakistan - Wikipedia](#)
https://en.wikipedia.org/wiki/East_Pakistan
 East Pakistan Wikipedia East Pakistan From Wikipedia the free encyclopedia Jump to navigation Jump to search Former province of Pakistan This article needs additional citations for verification P

[Pakistan Map / Geography of Pakistan / Map of Pakistan - Worldatlas.com](#)

User Interface and Comparisons with Google and Bing

Google Search

Google search results for "pakistan". The search bar shows "pakistan". Below it, there are tabs for All, News, Maps, Images, Videos, and More. The "All" tab is selected. The results page shows "About 1,220,000,000 results (0.97 seconds)". A "Top stories" section features three news items:

- In Pakistan's fight against Covid-19, religion might not be helping** (CNN.com, 1 day ago)
- COVID: 19:Pakistan fears rapid spread of coronavirus pandemic during Ramadan** (Gulf News, 12 hours ago)
- Covid-19: Pakistan's flag displayed at Switzerland's Matterhorn in solidarity** (The Express Tribune, 16 hours ago)

To the right of the stories is a box for "Pakistan" which includes the national flag, a map of South Asia, and a brief description: "Country in South Asia. Pakistan, officially the Islamic Republic of Pakistan, is a country in South Asia. It is the world's fifth-most populous country with a population exceeding 212.2 million. By area, it is the 33rd-largest country, spanning 881,913 square kilometres. [Wikimedia](#)".

Google Search

Google search results for "pakistan". The search bar shows "pakistan". Below it, there are tabs for All, News, Maps, Images, Videos, and More. The "All" tab is selected. The results page shows "www.britannica.com > place > Pakistan". A snippet from Britannica reads: "Pakistan | History - Geography | Britannica 3 days ago - Pakistan, populous and multiethnic country of South Asia. Having a predominately Indo-Iranian speaking population, Pakistan has historically ...".

A "Videos" section displays three video thumbnails:

- Pakistani Street Food Chicken Karahi Recipe!! | Street Food At ...** by Mark Wiens (23:46)
- When will Coronavirus Pandemic End in Pakistan? | Dr. Syed ...** by BOL News (1:25)
- LIVE BOL News Live | BOL News Live Streaming Pakistan ...** by BOL News (LIVE)

Google Search

Google search results for "pakistan". The search bar shows "pakistan". Below it, there are tabs for All, News, Maps, Images, Videos, and More. The "All" tab is selected. The results page shows "www.pakistan.gov.pk". A snippet from the site reads: "The Official Web Gateway to Pakistan LOGIN to e-Office | SIGN IN to my Mailbox ? bootstrap carousel - President of Pakistan - Prime Minister of Pakistan - National Assembly of Pakistan - Senate of ...".

Below this, there are links to "South Asia :: Pakistan – The World Factbook - Central ..." and "www.cnn.com · pakistan-coronavirus-lockdown-intl-hnk". A snippet from CNN reads: "In Pakistan's fight against Covid-19, religion might not be ... 2 days ago - Pakistan is entering its fifth week under lockdown in an effort to control the spread of the coronavirus."

An "Images for pakistan" section shows thumbnail images for various categories: flag, beautiful, Independence day, map, wallpaper, and 14 more.

User Interface and Comparisons with Google and Bing

Bing Search

pakistan

Sign in 200 Add the Give with Bing extension > Support nonprofits responding to COVID-19 when you search on Bing MAYBE LATER YES

What would you like to know about pakistan?

[pakistan news](#) [pakistan language](#) [pakistan time](#) [pakistan flag](#) [pakistan population](#)

The Official Web Gateway to Pakistan
pakistan.gov.pk ▾
A portal setup by Government of Pakistan to facilitate access to all the web sites run by various departments.
[The Official Web Gateway to Pakistan - Ministries](#)

News about Pakistan
bing.com/news

Pakistan's Anti-graft Body Issues Arrest Warrant Against Nawaz Sharif in Land

Bing Search

The Official Web Gateway to Pakistan
pakistan.gov.pk ▾
A portal setup by Government of Pakistan to facilitate access to all the web sites run by various departments.
[The Official Web Gateway to Pakistan - Ministries](#)

Pakistan - Wikipedia
https://en.wikipedia.org/Wiki/Pakistan ▾

[Overview](#) [Etymology](#) [History](#) [Role of Islam in Pakistan](#) [Geography, environment,](#) ▾
Pakistan, officially the Islamic Republic of Pakistan, is a country in South Asia. It is the world's fifth-most populous country with a population exceeding 212.2 million. By area, it is the 33rd-largest country, spanning 881,913 square kilometres (340,509 square miles). Pakistan has a 1,046-kilometre (650-mile) coastline along the Arabian Sea and Gulf of Oman in the south and is bordered by India to the east, Afghanistan to the west, Iran to the southwest, and China to the northeast. It is separated narrowly from Tajikistan by Afghanistan's Wakhan Corridor in the northwest, and als... +
Wikipedia · Text under CC-BY-SA license

Images of Pakistan
bing.com/images

Badshahi Mosque	Mohenjo-daro	Khyber Pass	Islamabad Capital Territory	Lahore Fort
Deosai National Park	Hingol National Park	Ayubia National Park	Khunjerab National Park	Kirthar National Park
Imran Khan Prime Minister	Muhammad Ali Jinnah	Benazir Bhutto	Arif Alvi President	Pervez Musharraf Former President

Data from: Wikipedia Cia Freebase
Wikimedia text under CC-BY-SA license
Suggest an edit

Bing Search

Images of Pakistan
bing.com/images

See more images of Pakistan

Pakistan | History - Geography | Britannica
https://www.britannica.com/Place/Pakistan ▾
Pakistan, populous and multilingual country of South Asia. Having a predominantly Indo-Iranian speaking population, Pakistan has historically and culturally been associated with ...

Pakistan travel | Asia - Lonely Planet
https://www.lonelyplanet.com/Pakistan ▾
Pakistan is the difficult child of South Asia – blessed with abundant natural and historical riches, but plagued by political instability, which has kept the country ...

5. Clustering (Anshul Pardhi)

Flat Clustering

K-means clustering algorithm is implemented for flat clustering of documents.

K-means clustering algorithm

Given a set of observations ($\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$), where each observation is a d -dimensional real vector, k -means clustering aims to partition the n observations into k ($\leq n$) sets $S = \{S_1, S_2, \dots, S_k\}$ so as to minimize the within-cluster sum of squares (WCSS). This WCSS is the sum of the squares of distances of all the observations belonging to a cluster from the centroid of the cluster. Formally, the objective is to find:

$$\arg \min_{\mathbf{S}} \sum_{i=1}^k \sum_{\mathbf{x} \in S_i} \|\mathbf{x} - \boldsymbol{\mu}_i\|^2$$

where $\boldsymbol{\mu}_i$ is the mean of points in S_i .

Selection of predefined clusters (value of k)

Values of k , ranging from 3 to 20 was tested on a sample of the dataset and the observed clusters were compared to the relevancy of the documents. Elbow method was also used to determine an appropriate value of k . After observing the results testing different values of k , a value of $k = 11$ gave most relevant results. As a result, all the fetched documents are stored in one of the 11 clusters.

Flat Clustering Algorithm Design & Implementation

All the documents crawled by Darrel (the student responsible for crawling) were indexed by Ruchi (the student responsible for indexing). Ruchi gave me the response generated by Apache Solr when a query "*" was given. This query retrieved all the crawled documents in a JSON format, and the results were dumped to a file. The size of this file was 1 GB.

- I extracted the "content" key's value from the JSON, which contained the content of the document in plain text format, removing the DOM and HTML tags.
- URL corresponding to the content was also extracted.
- The "content" value from all the crawled documents was stored in a list.
- The content was later vectorized using scikit-learn's TfidfVectorizer.

CLUSTERING

- Scikit-learn's KMeans was then used to apply the K-means clustering algorithm, with a value of k chosen as 11.
- The clusters obtained for a document after applying K-means algorithm as well as the URLs corresponding to that document were converted to Pandas series.
- The results were stored in a text file in the following format
 - URL, ClusterNumber

A sample of the generated clusters from the text file is shown below.

https://en.wikipedia.org/wiki/Foreign_relations_of_Afghanistan,1.0

https://en.wikipedia.org/wiki/Foreign_relations_of_Algeria,6.0

https://en.wikipedia.org/wiki/Foreign_relations_of_Antigua_and_Barbuda,4.0

https://en.wikipedia.org/wiki/Foreign_relations_of_Argentina,4.0

https://en.wikipedia.org/wiki/Foreign_relations_of_Armenia,7.0

https://en.wikipedia.org/wiki/Foreign_relations_of_Australia,8.0

https://en.wikipedia.org/wiki/Foreign_relations_of_Azerbaijan,1.0

https://en.wikipedia.org/wiki/Foreign_relations_of_Brazil,4.0

https://en.wikipedia.org/wiki/Foreign_relations_of_Denmark,7.0

https://en.wikipedia.org/wiki/Foreign_relations_of_Finland,7.0

https://en.wikipedia.org/wiki/Foreign_relations_of_France,7.0

https://en.wikipedia.org/wiki/Foreign_relations_of_Germany,7.0

https://en.wikipedia.org/wiki/Foreign_relations_of_Morocco,6.0

https://en.wikipedia.org/wiki/Foreign_relations_of_New_Zealand,8.0

https://en.wikipedia.org/wiki/Foreign_relations_of_Niue,8.0

https://en.wikipedia.org/wiki/Foreign_relations_of_Northern_Cyprus,1.0

https://en.wikipedia.org/wiki/Foreign_relations_of_Norway,7.0

https://en.wikipedia.org/wiki/Foreign_relations_of_Pakistan,1.0

https://en.wikipedia.org/wiki/Foreign_relations_of_Palestine,1.0

https://en.wikipedia.org/wiki/Foreign_relations_of_South_Korea,1.0

https://en.wikipedia.org/wiki/Foreign_relations_of_South_Sudan,6.0

https://en.wikipedia.org/wiki/Foreign_relations_of_the_Cook_Islands,8.0

CLUSTERING

As observed from the sample above, the clusters are created based on continents where a country belongs to (all Asian countries are grouped in cluster 1.0, all European countries are grouped in cluster 7.0 and so on).

The generated file contains cluster information of all the crawled pages.

Clustering results incorporation and integration with UI

- On the UI, there are 2 clustering modes: flat and hierarchical.
- If any of the clustering modes is selected, the GET request contains that clustering mode in the type parameter.
- Solr fetches the top 50 relevant results based on page ranking and HITS score and sends it to clustering computation handler.
- Inside the clustering computation handler, the results are ranked in such a way that cluster to which the top-ranked result belongs to, all the other results belonging to the same cluster in the order of decreasing page rank and HITS scores are ranked next.
- The process repeats for the next highest ranked result and all the pages belonging to its cluster are ranked next.
- This process repeats till we reach our threshold of 50.
- The updated ranks are then sent as a response to the UI.

Hierarchical Clustering

Agglomerative clustering algorithm with single link is implemented for hierarchical clustering of documents.

Agglomerative clustering algorithm

Hierarchical clustering is a method of cluster analysis which seeks to build a hierarchy of clusters. Agglomerative clustering is a bottom up approach of hierarchical clustering. Here, each observation, document here, starts in its own cluster, and pairs of clusters are merged as one moves up the hierarchy. The results of hierarchical clustering are usually represented by a dendrogram.

For computing distance, there are various distance metrics to choose from: Euclidean, Manhattan, Mahalanobis, to name a few. For the project, I have used Euclidean distance metric.

For computing the linkage criteria, to compute the distance between different clusters, there are a lot of metrics to choose from: maximum or complete link, minimum or single link, unweighted average link, weighted average link, centroid link, ward etc. For the project, I have used single link metric.

Hierarchical Clustering Algorithm Design & Implementation

All the documents crawled by Darrel (the student responsible for crawling) were indexed by Ruchi (the student responsible for indexing). Ruchi gave me the response generated by Apache Solr when a query “*” was given. This query retrieved all the crawled documents in a JSON format, and the results were dumped to a file. The size of this file was 1 GB.

- I extracted the “content” key’s value from the JSON, which contained the content of the document in plain text format, removing the DOM and HTML tags.
- URL corresponding to the content was also extracted.
- The “content” value from all the crawled documents was stored in a list.
- The content was later vectorized using scikit-learn’s TfidfVectorizer.
- Vectorizer’s output was converted to match agglomerative clustering algorithm’s input.
- Python’s fastcluster library is used to efficiently compute the results of hierarchical clustering faster.
- Firstly, linkage is calculated, and then the single method of fastcluster library is called to compute the clusters to which a URL belongs to.
- A dendrogram is plotted as well to show the results of clustering visually.
- The clusters obtained for a document after applying agglomerative clustering algorithm as well as the URLs corresponding to that document were converted to Pandas series.
- The results were stored in a text file in the following format
 - URL, ClusterNumber

A sample of the generated clusters from the text file is shown below.

<https://www.worldatlas.com/webimage/countrys/namerica/usstates/sc.htm>,6.0

<https://www.worldatlas.com/webimage/countrys/namerica/usstates/ustravel.htm>,6.0

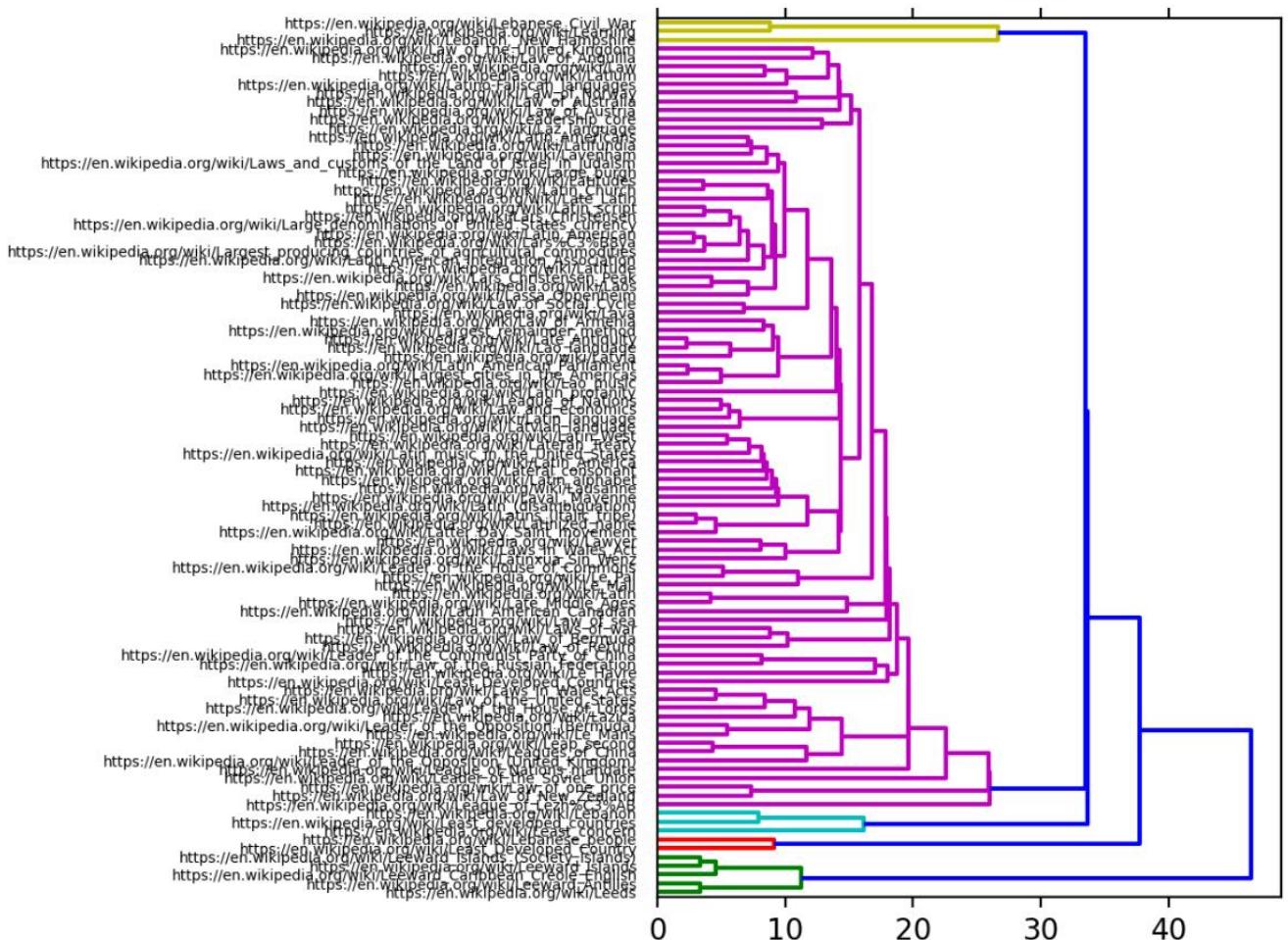
<https://www.worldatlas.com/webimage/countrys/oceania/mp.htm>,1.0

<https://www.worldatlas.com/webimage/countrys/oceania/fm.htm>,1.0

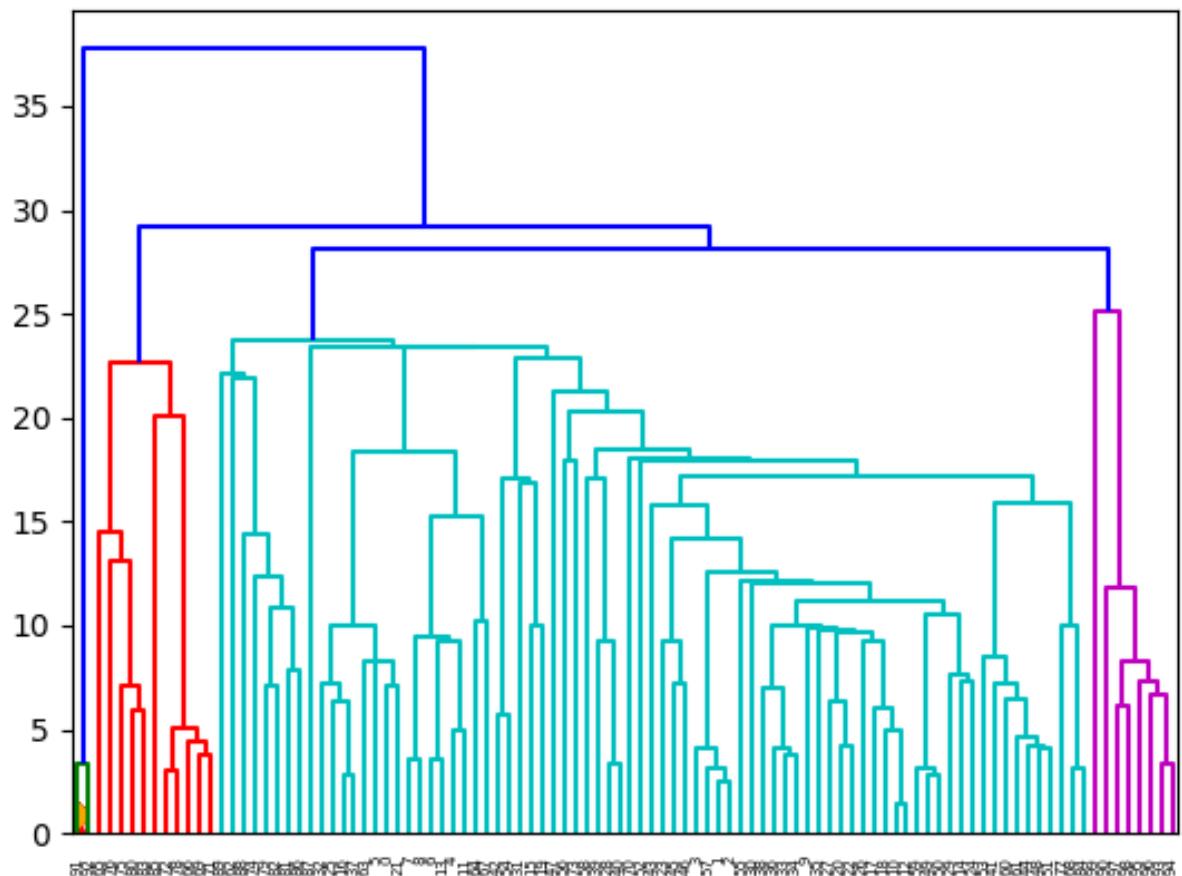
<https://www.worldatlas.com/webimage/countrys/oceania/tk.htm>,1.0

Since a dendrogram containing every links can get a lot cluttered and difficult to visualize, I have generated a couple of dendograms, taking samples from the data. A dendrogram containing all the webpages is shown as well in the end.

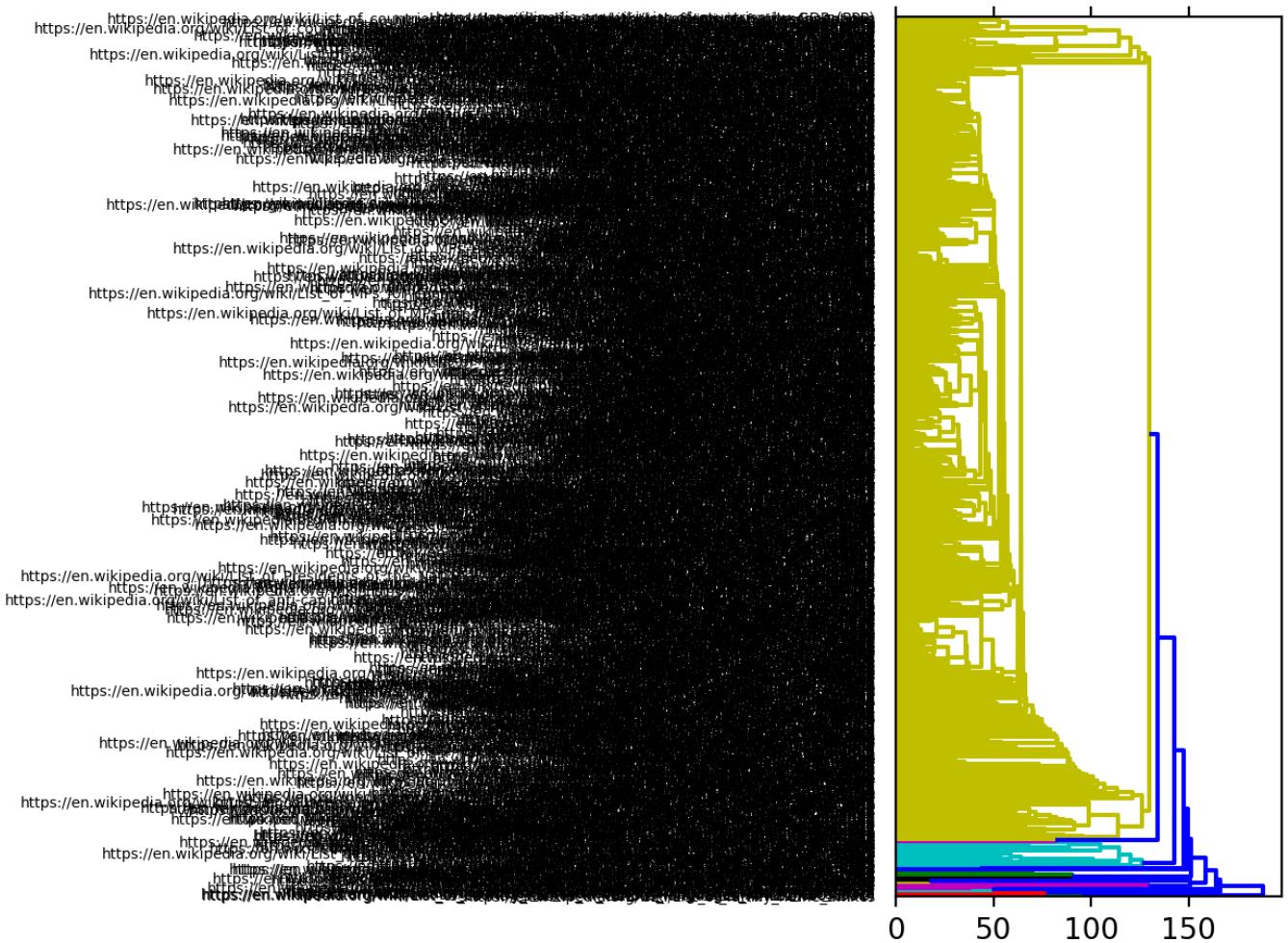
CLUSTERING



CLUSTERING



CLUSTERING



Number of clusters obtained

It was observed that fastcluster's agglomerative clustering selects the best number of clusters depending on the type of input data provided, and the various distance and linkage metrics used. It was observed that for our dataset, 7 clusters were created and all the links fall into one of these 7 clusters.

Clustering results incorporation and integration with UI

The steps are like those for flat clustering. There are two types of values for the "type" parameter in the request query from the UI, "flat_clustering" and "hierarchical_clustering". Depending on the clustering mode selected, that particular mode's clusters are selected for the re-computation of ranks of retrieved pages.

Number of queries experimented with to improve results

Around 50 queries were used per clustering method to improve the results obtained. Some of the results obtained after running these queries will be listed below in a further subsection. Observations are also provided below the screenshots which shows the impact of clustering on improving the relevancy of the retrieved results.

Query testing generation and impact on results and relevancy

I have used 50 queries for each of the two clustering methods: flat and hierarchical, to test the impact of the results of each clustering method. The queries were generated manually, depending on the initial results obtained from the original relevance model, using page rank and HITS. The results obtained after applying clustering seem relevant, as similar results are grouped together to be shown earlier than the non-relevant results.

Selection of queries

The criteria for selection of queries is based on the observation obtained from the type of clusters formed for both the clustering methods used. As shown in the sample subset of the results obtained above, one of the query selection criteria is continents, as it was observed that countries belonging to same continents were grouped together. Another selection criteria is the topic asked, such as map or constitution etc. For example, maps of different countries are grouped together while constitutions of different countries are grouped together.

Query examples with clustering

About 50 queries for both flat as well as hierarchical were tested. I have attached the results for 5 such queries below. The first image shows the results obtained by applying page rank. The second image shows the results obtained by applying flat clustering. The third image shows the results obtained after applying hierarchical clustering.

Query 1: Bangladesh constitution

CLUSTERING



Countries Search Engine

bangladesh constitution

Relevance Model Options:

Page Rank HITS

Clustering Options:

Flat Clustering Hierarchical Clustering

Query Expansion Option:

Association Metric Scalar

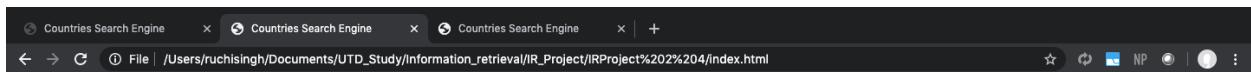
[Constitution of Bangladesh - Wikipedia](#)
https://en.wikipedia.org/wiki/Constitution_of_Bangladesh
Constitution of Bangladesh Wikipedia Constitution of Bangladesh From Wikipedia the free encyclopedia Jump to navigation Jump to search Constitution of the People's Republic of Bangladesh

[Socialism in Bangladesh - Wikipedia](#)
https://en.wikipedia.org/wiki/Socialism_in_Bangladesh
Socialism in Bangladesh Wikipedia Socialism in Bangladesh From Wikipedia the free encyclopedia Jump to navigation Jump to search Bangladesh This article is part of a series on the politics and government of Bangladesh

[President of Bangladesh - Wikipedia](#)
https://en.wikipedia.org/wiki/President_of_Bangladesh
President of Bangladesh Wikipedia President of Bangladesh From Wikipedia the free encyclopedia Jump to navigation Jump to search President of the People's Republic of Bangladesh

[Jatiya Sangsad - Wikipedia](#)
https://en.wikipedia.org/wiki/Jatiya_Sangsad
Jatiya Sangsad Wikipedia Jatiya Sangsad From Wikipedia the free encyclopedia Jump to navigation Jump to search Unicameral legislature of Bangladesh Sangsad redirects here For other uses see San

[Prime Minister of Bangladesh - Wikipedia](#)



Countries Search Engine

bangladesh constitution

Relevance Model Options:

Page Rank HITS

Clustering Options:

Flat Clustering Hierarchical Clustering

Query Expansion Option:

Association Metric Scalar

[Constitution of Bangladesh - Wikipedia](#)
https://en.wikipedia.org/wiki/Constitution_of_Bangladesh
Constitution of Bangladesh Wikipedia Constitution of Bangladesh From Wikipedia the free encyclopedia Jump to navigation Jump to search Constitution of the People's Republic of Bangladesh

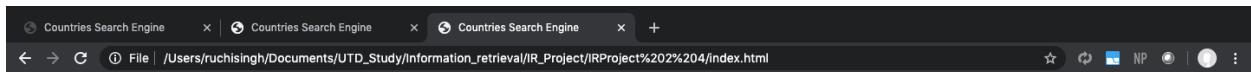
[President of Bangladesh - Wikipedia](#)
https://en.wikipedia.org/wiki/President_of_Bangladesh
President of Bangladesh Wikipedia President of Bangladesh From Wikipedia the free encyclopedia Jump to navigation Jump to search President of the People's Republic of Bangladesh

[Prime Minister of Bangladesh - Wikipedia](#)
https://en.wikipedia.org/wiki/Prime_Minister_of_Bangladesh
Prime Minister of Bangladesh Wikipedia Prime Minister of Bangladesh From Wikipedia the free encyclopedia Jump to navigation Jump to search Head of government of the People's Republic This article is about the current Prime Minister of Bangladesh

[Politics of Bangladesh - Wikipedia](#)
https://en.wikipedia.org/wiki/Politics_of_Bangladesh
Politics of Bangladesh Wikipedia Politics of Bangladesh From Wikipedia the free encyclopedia Jump to navigation Jump to search Political system of Bangladesh This article is about political system

[National Emblem of Bangladesh - Wikipedia](#)

CLUSTERING



Countries Search Engine

bangladesh constitution

Relevance Model Options:

Page Rank HITS

Clustering Options:

Flat Clustering Hierarchical Clustering

Query Expansion Option:

Association Metric Scalar

[Constitution of Bangladesh - Wikipedia](#)
https://en.wikipedia.org/wiki/Constitution_of_Bangladesh
Constitution of Bangladesh Wikipedia Constitution of Bangladesh From Wikipedia the free encyclopedia Jump to navigation Jump to search Constitution of the People's Republic of Bangladesh

[Socialism in Bangladesh - Wikipedia](#)
https://en.wikipedia.org/wiki/Socialism_in_Bangladesh
Socialism in Bangladesh Wikipedia Socialism in Bangladesh From Wikipedia the free encyclopedia Jump to navigation Jump to search Bangladesh This article is part of a series on the politics and government of Bangladesh

[National Emblem of Bangladesh - Wikipedia](#)
https://en.wikipedia.org/wiki/Emblem_of_Bangladesh
National Emblem of Bangladesh Wikipedia National Emblem of Bangladesh From Wikipedia the free encyclopedia Redirected from Emblem of Bangladesh Jump to navigation Jump to search National Emblem

[Victory day of Bangladesh - Wikipedia](#)
https://en.wikipedia.org/wiki/Victory_Day_of_Bangladesh
Victory day of Bangladesh Wikipedia Victory day of Bangladesh From Wikipedia the free encyclopedia Redirected from Victory Day of Bangladesh Jump to navigation Jump to search Victory Day

[President of Bangladesh - Wikipedia](#)

Observations: Clustering results show information relevant to constitution such as information about Prime Minister, President, Emblem, Politics of Bangladesh etc. Such results aren't shown by simply applying page rank.

Query 2: Australia currency

CLUSTERING



Countries Search Engine

australia currency

Relevance Model Options:

Page Rank HITS

Clustering Options:

Flat Clustering Hierarchical Clustering

Query Expansion Option:

Association Metric Scalar

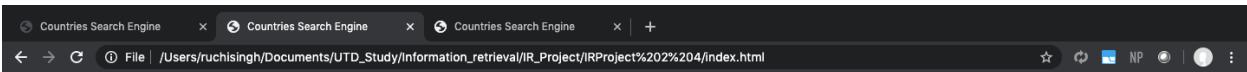
[Australian dollar - Wikipedia](#)
https://en.wikipedia.org/wiki/Australian_dollar
Australian dollar Wikipedia Australian dollar From Wikipedia the free encyclopedia Jump to navigation Jump to search Official currency used in Australia also used in Kiribati Nauru Tonga Tuvalu

[Local exchange trading system - Wikipedia](#)
https://en.wikipedia.org/wiki/Local_exchange_trading_system
Local exchange trading system Wikipedia Local exchange trading system From Wikipedia the free encyclopedia Jump to navigation Jump to search Local currency system This article needs additional cita

[What are the Most Expensive Currencies in the World? - Answers](#)
https://www.mapsofworld.com/answers/business/what-are-the-most-expensive-currencies-in-the-world/
What are the Most Expensive Currencies in the World Answers Maps of World Current Credible Consistent contactus mapsofworld com Facebook Twitter Google Pinterest Share Show Nav

[Asian Monetary Unit - Wikipedia](#)
https://en.wikipedia.org/wiki/Asian_Monetary_Unit
Asian Monetary Unit Wikipedia Asian Monetary Unit From Wikipedia the free encyclopedia Jump to navigation Jump to search This article needs additional citations for verification Please help impro

[Legal tender - Wikipedia](#)



Countries Search Engine

australia currency

Relevance Model Options:

Page Rank HITS

Clustering Options:

Flat Clustering Hierarchical Clustering

Query Expansion Option:

Association Metric Scalar

[Australian dollar - Wikipedia](#)
https://en.wikipedia.org/wiki/Australian_dollar
Australian dollar Wikipedia Australian dollar From Wikipedia the free encyclopedia Jump to navigation Jump to search Official currency used in Australia also used in Kiribati Nauru Tonga Tuvalu

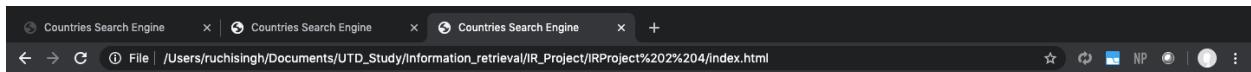
[Asian Monetary Unit - Wikipedia](#)
https://en.wikipedia.org/wiki/Asian_Monetary_Unit
Asian Monetary Unit Wikipedia Asian Monetary Unit From Wikipedia the free encyclopedia Jump to navigation Jump to search This article needs additional citations for verification Please help impro

[Legal tender - Wikipedia](#)
https://en.wikipedia.org/wiki/Legal_tender
Legal tender Wikipedia Legal tender From Wikipedia the free encyclopedia Jump to navigation Jump to search This article is about the payment medium For the song see Legal Tender song This art

[Currency union - Wikipedia](#)
https://en.wikipedia.org/wiki/Monetary_union
Currency union Wikipedia Currency union From Wikipedia the free encyclopedia Redirected from Monetary union Jump to navigation Jump to search Part of a series on World trade Policy Import Expo

[Tuvaluan dollar - Wikipedia](#)

CLUSTERING



Countries Search Engine

australia currency

Relevance Model Options:

Page Rank HITS

Clustering Options:

Flat Clustering Hierarchical Clustering

Query Expansion Option:

Association Metric Scalar

[Australian dollar - Wikipedia](#)
https://en.wikipedia.org/wiki/Australian_dollar
Australian dollar Wikipedia Australian dollar From Wikipedia the free encyclopedia Jump to navigation Jump to search Official currency used in Australia also used in Kiribati Nauru Tonga Tuvalu

[Asian Monetary Unit - Wikipedia](#)
https://en.wikipedia.org/wiki/Asian_Monetary_Unit
Asian Monetary Unit Wikipedia Asian Monetary Unit From Wikipedia the free encyclopedia Jump to navigation Jump to search This article needs additional citations for verification Please help impro

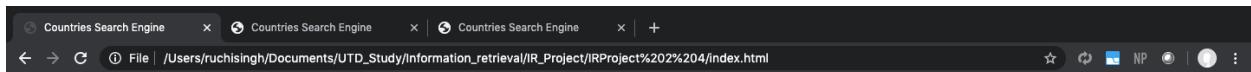
[Counterfeit money - Wikipedia](#)
https://en.wikipedia.org/wiki/Counterfeit_money
Counterfeit money Wikipedia Counterfeit money From Wikipedia the free encyclopedia Jump to navigation Jump to search See also Coin counterfeiting and Slug coin This article has multiple issues

[ASEAN Free Trade Area - Wikipedia](#)
https://en.wikipedia.org/wiki/ASEAN_Free_Trade_Area
ASEAN Free Trade Area Wikipedia ASEAN Free Trade Area From Wikipedia the free encyclopedia Jump to navigation Jump to search Free trade area of the Association of South East Asian Nations This art

Observations: Clustering results show results like currency such as monetary unit, legal tender, currency union, counterfeit money etc. Furthermore, these results show similar results of countries near Australia such as Tuvalu.

Query 3: European Union

CLUSTERING



Countries Search Engine

europen union

Relevance Model Options:

Page Rank HITS

Clustering Options:

Flat Clustering Hierarchical Clustering

Query Expansion Option:

Association Metric Scalar

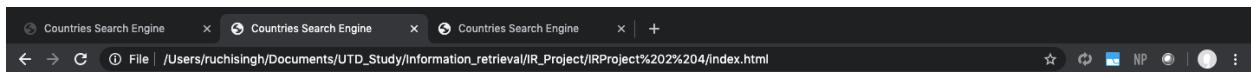
[List of countries by leading trade partners - Wikipedia](#)
https://en.wikipedia.org/wiki/List_of_countries_by_leading_trade_partners
List of countries by leading trade partners Wikipedia List of countries by leading trade partners From Wikipedia the free encyclopedia Jump to navigation Jump to search Wikipedia list article Part

[Portal:European Union - Wikipedia](#)
https://en.wikipedia.org/wiki/Portal:European_Union
Portal European Union Wikipedia Portal European Union From Wikipedia the free encyclopedia Jump to navigation Jump to search Wikimedia portal Portal maintenance status February This portal

[European Union Customs Union - Wikipedia](#)
https://en.wikipedia.org/wiki/European_Union_Customs_Union
European Union Customs Union Wikipedia European Union Customs Union From Wikipedia the free encyclopedia Jump to navigation Jump to search This article is about European Union Customs Union It is

[Treaties of the European Union - Wikipedia](#)
https://en.wikipedia.org/wiki/Treaties_of_the_European_Union
Treaties of the European Union Wikipedia Treaties of the European Union From Wikipedia the free encyclopedia Jump to navigation Jump to search Treaties of the European Union Front page of an EU doc

[Institutions of the European Union - Wikipedia](#)



Countries Search Engine

europen union

Relevance Model Options:

Page Rank HITS

Clustering Options:

Flat Clustering Hierarchical Clustering

Query Expansion Option:

Association Metric Scalar

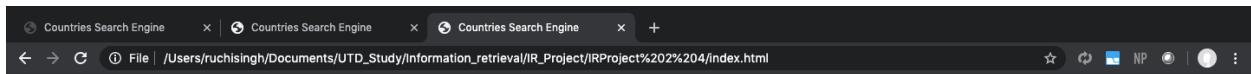
[List of countries by leading trade partners - Wikipedia](#)
https://en.wikipedia.org/wiki/List_of_countries_by_leading_trade_partners
List of countries by leading trade partners Wikipedia List of countries by leading trade partners From Wikipedia the free encyclopedia Jump to navigation Jump to search Wikipedia list article Part

[Withdrawal from the European Union - Wikipedia](#)
https://en.wikipedia.org/wiki/Withdrawal_from_the_European_Union
Withdrawal from the European Union Wikipedia Withdrawal from the European Union From Wikipedia the free encyclopedia Jump to navigation Jump to search Article redirects here For the former Eu

[Abbreviations of the European Union](#)
https://abbreviations.yourdictionary.com/articles/abbreviations-european-union.html
Abbreviations of the European Union This website uses cookies to ensure you get the best experience Learn more Got it Reference Menu Dictionary Thesaurus Examples Sentences Quotes Reference Spanish

[Police and Judicial Co-operation in Criminal Matters - Wikipedia](#)
https://en.wikipedia.org/wiki/Police_and_Judicial_Co-operation_in_Criminal_Matters
Police and Judicial Co operation in Criminal Matters Wikipedia Police and Judicial Co operation in Criminal Matters From Wikipedia the free encyclopedia Jump to navigation Jump to search PJCC red

CLUSTERING



Countries Search Engine

european union

Relevance Model Options:
 Page Rank HITS

Clustering Options:
 Flat Clustering Hierarchical Clustering

Query Expansion Option:
 Association Metric Scalar

[List of countries by leading trade partners - Wikipedia](https://en.wikipedia.org/wiki/List_of_countries_by_leading_trade_partners)
https://en.wikipedia.org/wiki/List_of_countries_by_leading_trade_partners
List of countries by leading trade partners Wikipedia List of countries by leading trade partners From Wikipedia the free encyclopedia Jump to navigation Jump to search Wikipedia list article Part

[Abbreviations of the European Union](https://en.wikipedia.org/wiki/Abbreviations_of_the_European_Union)
https://abbreviations.yourdictionary.com/articles/abbreviations-european-union.html
Abbreviations of the European Union This website uses cookies to ensure you get the best experience Learn more Got it Reference Menu Dictionary Thesaurus Examples Sentences Quotes Reference Spanish

[Maastricht Treaty - Wikipedia](https://en.wikipedia.org/wiki/Maastricht_Treaty)
https://en.wikipedia.org/wiki/Maastricht_Treaty
Maastricht Treaty Wikipedia Maastricht Treaty From Wikipedia the free encyclopedia Jump to navigation Jump to search Founding treaty of the European Union This article is about one of two founding

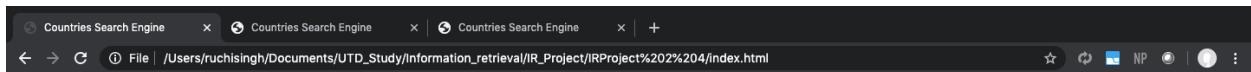
[Portal:European Union - Wikipedia](https://en.wikipedia.org/wiki/Portal:European_Union)
https://en.wikipedia.org/wiki/Portal:European_Union
Portal European Union Wikipedia Portal European Union From Wikipedia the free encyclopedia Jump to navigation Jump to search Wikimedia portal Portal maintenance status February This portal

[Treaties of the European Union - Wikipedia](https://en.wikipedia.org/wiki/Treaties_of_the_European_Union)
https://en.wikipedia.org/wiki/Treaties_of_the_European_Union

Observations: Hierarchical clustering shows a result Maastricht Treaty. It was the treaty that led to the creation of the European Union. Such results are missing from the top results of page rank.

Query 4: Mumbai

CLUSTERING



Countries Search Engine

mumbai

Relevance Model Options:

Page Rank HITS

Clustering Options:

Flat Clustering Hierarchical Clustering

Query Expansion Option:

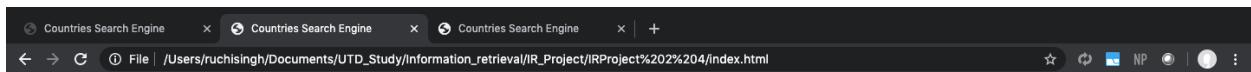
Association Metric Scalar

Mumbai - Wikipedia
https://en.wikipedia.org/wiki/Mumbai
Mumbai Wikipedia Mumbai From Wikipedia the free encyclopedia Jump to navigation Jump to search Bombay redirects here For other uses see Bombay disambiguation and Mumbai disambiguation Meg

Mumbai Population 2020 (Demographics, Maps, Graphs)
https://worldpopulationreview.com/world-cities/mumbai-population/
Mumbai Population Demographics Maps Graphs World Population Review Mobile Navigation Home Continents Countries World Cities US States US Counties US Cities Zips Canadian Provinces Coronavirus

Administrative divisions of Mumbai - Wikipedia
https://en.wikipedia.org/wiki/Administrative_divisions_of_Mumbai
Administrative divisions of Mumbai Wikipedia Administrative divisions of Mumbai From Wikipedia the free encyclopedia Jump to navigation Jump to search This article needs additional citations for ve

History of Mumbai from an Arabian seaport to Bollywood | Britannica
https://www.britannica.com/video/141204/Video-overview-Mumbai-history
History of Mumbai from an Arabian seaport to Bollywood Britannica Search Britannica Encyclop dia Britannica Login Categories Science Technology Health Medicine Sports Recreation Geography Trav



Countries Search Engine

mumbai

Relevance Model Options:

Page Rank HITS

Clustering Options:

Flat Clustering Hierarchical Clustering

Query Expansion Option:

Association Metric Scalar

Mumbai - Wikipedia
https://en.wikipedia.org/wiki/Mumbai
Mumbai Wikipedia Mumbai From Wikipedia the free encyclopedia Jump to navigation Jump to search Bombay redirects here For other uses see Bombay disambiguation and Mumbai disambiguation Meg

Mumbai Population 2020 (Demographics, Maps, Graphs)
https://worldpopulationreview.com/world-cities/mumbai-population/
Mumbai Population Demographics Maps Graphs World Population Review Mobile Navigation Home Continents Countries World Cities US States US Counties US Cities Zips Canadian Provinces Coronavirus

The 10 Largest Cities in India - WorldAtlas.com
https://www.worldatlas.com/articles/the-10-largest-cities-in-india.html
The Largest Cities in India WorldAtlas com Trending Covid Cases And Deaths Per State In The US Could There Be Another Great Depression Celebrities Who Have Recovered From COVID Society Th

The 10 Largest Cities in the World - WorldAtlas.com
https://www.worldatlas.com/amp/articles/the-10-largest-cities-in-the-world.html
The Largest Cities in the World WorldAtlas com Trending Covid Cases And Deaths Per State In The US Could There Be Another Great Depression Celebrities Who Have Recovered From COVID World

CLUSTERING

The screenshot shows a web browser window titled "Countries Search Engine" with three tabs open, all showing the same search results page. The URL in the address bar is "/Users/ruchisingh/Documents/UTD_Study/Information_retrieval/IR_Project/IRProject%202%204/index.html". The search query "mumbai" is entered in the search bar. The results are displayed in a table with columns for "Result", "URL", and "Description".

Result	URL	Description
Mumbai - Wikipedia	https://en.wikipedia.org/wiki/Mumbai	Mumbai Wikipedia Mumbai From Wikipedia the free encyclopedia Jump to navigation Jump to search Bombay redirects here For other uses see Bombay disambiguation and Mumbai disambiguation Meg
Mumbai Population 2020 (Demographics, Maps, Graphs)	https://worldpopulationreview.com/world-cities/mumbai-population/	Mumbai Population Demographics Maps Graphs World Population Review Mobile Navigation Home Continents Countries World Cities US States US Counties US Cities Zips Canadian Provinces Coronavirus
Bene Israel - Wikipedia	https://en.wikipedia.org/wiki/Bene_Israel	Bene Israel Wikipedia Bene Israel From Wikipedia the free encyclopedia Jump to navigation Jump to search Jewish ethnic group Not to be confused with Bani Isra il or Beta Israel Bene Israel
India geography, maps, climate, environment and terrain from India CountryReports - CountryReports	https://www.countryreports.org/country/India/geography.htm	India geography maps climate environment and terrain from India CountryReports CountryReports India Overview People Government Politics Geography Maps Geography Comparisons Economy News

Observations: Clustering not only improves results for countries, but it also improves results for cities. Here flat clustering groups relevant information like the top 10 largest cities in India as well as the world. Furthermore, hierarchical clustering shows a result, Bene Israel. These people are the Jewish community settled predominantly in Mumbai. Such results aren't ranked that high when applying page rank.

Query 5: Guatemala

CLUSTERING



Countries Search Engine

guatemala

Relevance Model Options:

Page Rank HITS

Clustering Options:

Flat Clustering Hierarchical Clustering

Query Expansion Option:

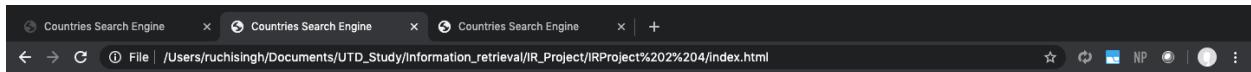
Association Metric Scalar

[Guatemala Map | Map of Guatemala](#)
https://www.mapsofworld.com/guatemala/
Guatemala Map Map of Guatemala Menu World Map World Maps Political Map of the World Physical Map of the World Blank World Map World Map for Kids Earth Map World Atlas World Time Zones Map Coun

[Guatemala City - Wikipedia](#)
https://en.wikipedia.org/wiki/Guatemala_City
Guatemala City Wikipedia Guatemala City From Wikipedia the free encyclopedia Jump to navigation Jump to search Capital of Guatemala Capital City in Guatemala Guatemala Guatemala City Guatemala Cap

[Guatemala Population 2020 \(Demographics, Maps, Graphs\)](#)
https://worldpopulationreview.com/countries/guatemala-population/
Guatemala Population Demographics Maps Graphs World Population Review Mobile Navigation Home Continents Countries World Cities US States US Counties US Cities Zips Canadian Provinces Coronavi

[Guatemala - Wikipedia](#)
https://en.wikipedia.org/wiki/Guatemala
Guatemala Wikipedia Guatemala From Wikipedia the free encyclopedia This is the latest accepted revision reviewed on April Jump to navigation Jump to search This article is about the cou



Countries Search Engine

guatemala

Relevance Model Options:

Page Rank HITS

Clustering Options:

Flat Clustering Hierarchical Clustering

Query Expansion Option:

Association Metric Scalar

[Guatemala Map | Map of Guatemala](#)
https://www.mapsofworld.com/guatemala/
Guatemala Map Map of Guatemala Menu World Map World Maps Political Map of the World Physical Map of the World Blank World Map World Map for Kids Earth Map World Atlas World Time Zones Map Coun

[Current Local Time in Guatemala City, Guatemala](#)
https://www.timeanddate.com/worldclock/guatemala/guatemala
Current Local Time in Guatemala City Guatemala Menu timeanddate.com Search Site Articles City Country Social Share this page Follow us on Facebook Tweet Follow Twitter Visit us on Faceboo

[Mayan | Ethnologue](#)
https://www.ethnologue.com/subgroups/mayan
Mayan Ethnologue Skip to main content Login Shopping cart Ethnologue Languages Countries Guides About Plans Pricing Mayan Print Mayan Huastecan Chicomecilete cob A language of Mexico

[Guatemala City Population 2020 \(Demographics, Maps, Graphs\)](#)
https://worldpopulationreview.com/world-cities/guatemala-city-population/
Guatemala City Population Demographics Maps Graphs World Population Review Mobile Navigation Home Continents Countries World Cities US States US Counties US Cities Zips Canadian Provinces Cor

CLUSTERING

guatemala

Relevance Model Options:

Page Rank HITS

Clustering Options:

Flat Clustering Hierarchical Clustering

Query Expansion Option:

Association Metric Scalar

Search

[Guatemala Map | Map of Guatemala](#)
https://www.mapsofworld.com/guatemala/
Guatemala Map Map of Guatemala Menu World Map World Maps Political Map of the World Physical Map of the World Blank World Map World Map for Kids Earth Map World Atlas World Time Zones Map Coun

[Mayan | Ethnologue](#)
https://www.ethnologue.com/subgroups/mayan
Mayan Ethnologue Skip to main content Login Shopping cart Ethnologue Languages Countries Guides About Plans Pricing Mayan Print Mayan Huastecan Chicomuceltec cob A language of Mexico

[Guatemala | History, Map, Flag, Population, & Facts | Britannica](#)
https://www.britannica.com/place/Guatemala
Guatemala History Map Flag Population Facts Britannica Search Britannica Encyclop dia Britannica Login Categories Science Technology Health Medicine Sports Recreation Geography Travel

[Guatemala Internet Usage and Telecom Statistics](#)
https://www.internetworldstats.com/am/gt.htm
Guatemala Internet Usage and Telecom Statistics Central America Internet Stats Site Links Guatemala Guatemala Internet usage broadband and telecommunications reports Internet Usage Statisti

[Guatemala City - Wikipedia](#)

Observations: The majority of indigenous people in Guatemala are of Mayan origin. Clustering handles such details and shows Mayan results on querying Guatemala. Such results don't get a top rank when page rank is used.

6.Query Expansion and Relevance Feedback

Query expansion is used to improve search results by following local methods:

- 1) Relevance Feedback (Rocchio Algorithm)
- 2) Pseudo Relevance Feedback
 - a) Association cluster
 - b) Metric cluster
 - c) Scalar cluster

The following formula is used to implement Rocchio algorithm:

- Used in practice:

$$\vec{q}_m = \alpha \vec{q}_0 + \beta \frac{1}{|D_r|} \sum_{\vec{d}_j \in D_r} \vec{d}_j - \gamma \frac{1}{|D_{nr}|} \sum_{\vec{d}_j \in D_{nr}} \vec{d}_j$$

- D_r = set of known relevant doc vectors
- D_{nr} = set of known irrelevant doc vectors
 - Different from C_r and C_{nr} 
- \vec{q}_m = modified query vector; \vec{q}_0 = original query vector; α, β, γ : weights (hand-chosen or set empirically)
- New query moves toward relevant documents and away from irrelevant documents

According to this image, 20 queries were fed as q_0 as input and the algorithm was tested for alpha =1.0, beta = 0.9 and gamma =0.1 values , so that a set of weights could be found that can separate relevant documents from the irrelevant documents.

The 20 queries were selected based on number of query parameters, spelling mistakes, regional name of countries and queries containing words that are states or cities (outside of country domain).

- 1) 4 queries with one keyword for the data that was crawled
 - a) australia
 - b) china
 - c) ethiopia
 - d) pakistan
- 2) 4 queries with two keywords for the data that was crawled
 - a) bangladesh area
 - b) pakistan terrorism
 - c) china population
 - d) india pakistan
- 3) 4 queries with spelling mistakes in the query words for the data that was crawled
 - a) australai
 - b) paksitan
 - c) chna
 - d) germny
- 4) 4 queries with regional name for countries
 - a) hindustan
 - b) britaan
 - c) franche

- d) gana
- 5) 4 queries with the words that are outside of country domain
- texas
 - europe
 - mumbai
 - doha

examples of the web pages that are found relevant:

Query	Relevant web pages
australia	https://en.wikipedia.org/wiki/Australia
	https://www.britannica.com/place/Australia
	https://en.wikipedia.org/wiki/Demographics_of_Australia
	https://en.wikipedia.org/wiki/Indigenous_Australian_music
china population	https://worldpopulationreview.com/countries/china-population/
	https://en.wikipedia.org/wiki/Chinese_people
	https://en.wikipedia.org/wiki/One-child_policy
	https://www.statista.com/topics/753/china/
india pakistan	https://en.wikipedia.org/wiki/Pakistan_Movement
	https://en.wikipedia.org/wiki/Kashmir_conflict
	https://worldview.stratfor.com/topic/india-pakistan-rivalry

examples of the web pages that are found irrelevant:

Query	Irrelevant web pages
australia	https://en.wikipedia.org/wiki/Cocos_Islands
	https://en.wikipedia.org/wiki/Demographics_of_Cocos_(Keeling)_Islands
	https://en.wikipedia.org/wiki/Applied_linguistics
china population	https://en.wikipedia.org/wiki/Taiwan,_China
	http://www.spacetoday.org/China/ChinaGlance.html
	https://www.thoughtco.com/peoples-republic-of-china-facts-history-195233
india pakistan	https://en.wikipedia.org/wiki/Victory_Day_of_Bangladesh
	https://en.wikipedia.org/wiki/Biharis_in_Bangladesh
	https://en.wikipedia.org/wiki/Gilgit-Baltistan

modified queries after applying Roccio algorithm to above queries:

Original query	Query after roccio algorithm
australia	australia map
china	china republic
ethiopia	ethiopia map
pakistan	pakistan kashmir
bangladesh area	bangladesh area code
pakistan terrorism	pakistan terrorism news
china population	china population policy
india pakistan	india Pakistan bbc
australai	australai coco
paksitan	paksitan east
chna	chna corona

germny	germny map
hindustan	hindustan time
britaan	britaan bhutan
franche	franche fries
gana	gana gana
indo	indo democrat
indo china	indo china wall
image argentina	imagining argentina wiki
iraq for sale	iraq sale people

The approach was used for the 20 queries mentioned above. The results obtained were not very promising and universal set of weights could not be arrived. Following were some of the findings:

1. Some queries fetched no documents – like “hindustan”, “britaan” etc. which meant that the data we crawled did not contain neither the pages nor any reference to these words. In such cases, relevance model did not work
2. Some queries fetched all irrelevant documents – like “image argentina” fetched results for “imagining argentina movie”. So the modified query was “image argentina satelite”
3. Some queries fetched great relevant results – like “china population” which fetched the china population 2020 and china child policy. In such cases, weights stabilized too early and so the modified query was not much relevant.
4. Performance issues as too much time was taken for some queries to get the stabilized weights, hence could not be performed on the fly
5. Each query needed its own unique set of weights

Due to the above issues, it was decided not to include the code for Rocchio algorithm for expanding the query, rather the results of the observations were used in the pseudo relevance feedback approaches to filter out the stop words, avoid spelling mistakes etc.

Pseudo Relevance Feedback:

1. **Association cluster:** The idea is that stems which co-occur frequently inside documents have a synonymity association
2. **Metric cluster:** The idea is that the stems which occur far apart in the document are less co-related to the terms that occur closer like in the same sentence.
3. **Scalar cluster:** The idea is that two stems are more co-related if they have similar neighbourhood.

Following is the list of 54 queries used for pseudo relevance feedback. (Each query run 3 times with 3 different cluster methods as mentioned above):

Query	Relevant results fetched after expansion(pages)	Most relevant expanded query	Cluster method used
china	10	china map republic	association, scalar
china population	10	china population policy child	association,
australia	10	australia coco demographic	association, scalar
bangladesh	9	bangladesh people arm	metric, association
ethiopia	2	ethiopia concur confluence	metric, association, scalar

india pakistan	8	india pakistan kashmir bangladesh war	metric, association
france	5	france wins wine map	metric, association, scalar
germany	5	germany squad language human	metric, association, scalar
france area	8	france area code people	association, scalar
germany flag	10	germany flag color red	association
china coronavirus	9	china coronavirus spread world	metric, association, scalar
pakistan terrorist	4	pakistan terrorist news kashmir	metric, association, scalar
qatar	9	qatar airway oil currency	metric, association,
canada weather	1	canada weather daily toronto canadian	metric, association, scalar
india	9	india map democracy	association
hindustan	0	hindustan article time	metric, association, scalar
gana	0	gana gana	metric, association
franche	0	franche firewatch	association, scalar

Based upon the above results, it was decided that the association cluster method is the best for implementation of the pseudo relevance feedback.

Some challenges:

1. Scalar clustering took too long to run so we can't perform them well
2. Metric clustering also was taking time in some cases.

Query 1: australia

local document set(ids):

```
491d54c4-885c-11ea-bc55-0242ac130003
491d571c-885c-11ea-bc55-0242ac130003
491d57e4-885c-11ea-bc55-0242ac130003
491d58a2-885c-11ea-bc55-0242ac130003
491d5ab4-885c-11ea-bc55-0242ac130003
491d5b90-885c-11ea-bc55-0242ac130003
491d5c58-885c-11ea-bc55-0242ac130003
491d5d16-885c-11ea-bc55-0242ac130003
491d5f00-885c-11ea-bc55-0242ac130003
491d5fd2-885c-11ea-bc55-0242ac130003
491d6090-885c-11ea-bc55-0242ac130003
491d614e-885c-11ea-bc55-0242ac130003
491d620c-885c-11ea-bc55-0242ac130003
491d62ca-885c-11ea-bc55-0242ac130003
491d6392-885c-11ea-bc55-0242ac130003
491d65cc-885c-11ea-bc55-0242ac130003
491d66a8-885c-11ea-bc55-0242ac130003
491d6766-885c-11ea-bc55-0242ac130003
491d682e-885c-11ea-bc55-0242ac130003
491d68ec-885c-11ea-bc55-0242ac130003
491d69b4-885c-11ea-bc55-0242ac130003
```

491d6a72-885c-11ea-bc55-0242ac130003
491d6b30-885c-11ea-bc55-0242ac130003
491d6d6a-885c-11ea-bc55-0242ac130003
491d6e32-885c-11ea-bc55-0242ac130003
491d6ef0-885c-11ea-bc55-0242ac130003
491d6fae-885c-11ea-bc55-0242ac130003
491d7076-885c-11ea-bc55-0242ac130003
491d71ac-885c-11ea-bc55-0242ac130003
491d7288-885c-11ea-bc55-0242ac130003



india_pakistan_local_doc.txt

local vocabulary set:

local_vocabulary_set1.txt

local stem set:

local_stem_set1.txt

Query 2: india Pakistan**local document set(ids):**

867c15ba-885e-11ea-bc55-0242ac130003
867c193e-885e-11ea-bc55-0242ac130003
867c1b0a-885e-11ea-bc55-0242ac130003
867c1c72-885e-11ea-bc55-0242ac130003
867c1db2-885e-11ea-bc55-0242ac130003
867c1eca-885e-11ea-bc55-0242ac130003
867c205a-885e-11ea-bc55-0242ac130003
867c219a-885e-11ea-bc55-0242ac130003
867c229e-885e-11ea-bc55-0242ac130003
867c23fc-885e-11ea-bc55-0242ac130003
867c253c-885e-11ea-bc55-0242ac130003
867c264a-885e-11ea-bc55-0242ac130003
867c27d0-885e-11ea-bc55-0242ac130003
867c2906-885e-11ea-bc55-0242ac130003
867c2a1e-885e-11ea-bc55-0242ac130003
867c2b7c-885e-11ea-bc55-0242ac130003
867c2cc6-885e-11ea-bc55-0242ac130003
867c2e60-885e-11ea-bc55-0242ac130003
867c2faa-885e-11ea-bc55-0242ac130003
867c30fe-885e-11ea-bc55-0242ac130003
867c3252-885e-11ea-bc55-0242ac130003
867c3392-885e-11ea-bc55-0242ac130003
867c3518-885e-11ea-bc55-0242ac130003
867c36e4-885e-11ea-bc55-0242ac130003

867c38b0-885e-11ea-bc55-0242ac130003
867c39fa-885e-11ea-bc55-0242ac130003
867c3b3a-885e-11ea-bc55-0242ac130003
867c3c98-885e-11ea-bc55-0242ac130003
867c4350-885e-11ea-bc55-0242ac130003
867cac96-885e-11ea-bc55-0242ac130003



india_pakistan_local_doc.txt

local vocabulary set:



local_vocabulary_set2.txt

local stem set:



local_stem_set2.txt

Query 3: china population

local documents (ids) :

e5200812-8861-11ea-bc55-0242ac130003
e5200a10-8861-11ea-bc55-0242ac130003
e5200b0a-8861-11ea-bc55-0242ac130003
e520115e-8861-11ea-bc55-0242ac130003
e52014e2-8861-11ea-bc55-0242ac130003
e520160e-8861-11ea-bc55-0242ac130003
e52017d0-8861-11ea-bc55-0242ac130003
e5202130-8861-11ea-bc55-0242ac130003
e520241e-8861-11ea-bc55-0242ac130003
e520250e-8861-11ea-bc55-0242ac130003
e52025e0-8861-11ea-bc55-0242ac130003
e520269e-8861-11ea-bc55-0242ac130003
e520275c-8861-11ea-bc55-0242ac130003
e520281a-8861-11ea-bc55-0242ac130003
e52029b4-8861-11ea-bc55-0242ac130003
e5202a90-8861-11ea-bc55-0242ac130003
e5202b62-8861-11ea-bc55-0242ac130003
e5202c20-8861-11ea-bc55-0242ac130003
e5202e78-8861-11ea-bc55-0242ac130003
e5202f54-8861-11ea-bc55-0242ac130003
e520301c-8861-11ea-bc55-0242ac130003
e52030e4-8861-11ea-bc55-0242ac130003
e52031ac-8861-11ea-bc55-0242ac130003
e520326a-8861-11ea-bc55-0242ac130003
e5203328-8861-11ea-bc55-0242ac130003
e520362a-8861-11ea-bc55-0242ac130003
e5203710-8861-11ea-bc55-0242ac130003

e52037ce-8861-11ea-bc55-0242ac130003
e5203896-8861-11ea-bc55-0242ac130003
e520395e-8861-11ea-bc55-0242ac130003



china_population_local_doc.txt

local vocabulary set:

local_vocabulary_set3.txt

local stem set:

local_stem_set3.txt

from the table above, we can see that the relevant pages obtained are more when association cluster method is used and the time consumed for association cluster is less. As a result, we can say that association cluster is best suitable when performance is considered. It also takes care of spelling mistakes. The number of relevant pages obtained from table shows the correlations for the queries. It also shows the expanded query for each original query.

Collaboration with UI and relevance model:

- 1) From the API (URL) call we receive original query and the cluster method type to be used.
- 2) Based on these parameters, the relevant model provides top 50 documents to the query expansion python code. These results are then used as local documents to form local vocabulary , stem sets.
- 3) Based on the query parameter the cluster method out of association, metric and scalar is called to return expanded query to relevance model.
- 4) The relevance model then hits solar in order to received results for expanded query.
- 5) These results are then passed to front end and displayed.

Query selection for demo:

For the demo, the query “china population” was used since it has two words and can be enhanced with other words. It is tested for association cluster method.

7. Discussion

After completing this project, we have learned the many different aspects required to make a simple and working search engine. We were able to create a search engine that could take crawled webpages and return somewhat relevant results using the different techniques that we learned in class. The project gave us lots of practical knowledge of creating a search engine and also showed just that we have only really scratched the surface of Information Retrieval. When comparing the results that we worked so hard to get with large search engines such as Google and Bing, we realized that the work we did is small compared to what it takes to create and maintain popular and up-to-date search engines like Google or Bing.

8. Conclusion

In conclusion, for this project, we created a search engine — using the knowledge we gained throughout this class and some open source programs — that focuses on information related to Countries. The five of us splits the 5 different tasks required to build the search engine and collaborated to combine our individual parts into a final product that can provide relatively accurate information about Countries. While there were many challenges along the way, we were able to overcome them and create a good search engine. This project provided us with an opportunity to apply the knowledge we gained in this class in a practical way and boosted our understanding of the theoretical concepts we learned.