

The University of Texas at Dallas

***CS 6322 Information Retrieval Spring
2024***

Class Project Proposal

Project TITLE: Search Engine for Skiing

Group: No 5

Students:

*Aishwarya Vinod Menon,
axv220062@utdallas.edu*

*Anusha Gupta,
axg230026@utdallas.edu*

*Hanisha Anil Mohinani,
hxm220089@utdallas.edu*

*Sowmya Sivaramakrishnan,
sxs230043@utdallas.edu*

*Vidya Sreekumar,
vxs220066@utdallas.edu*

1. The Problem

To develop a search engine for skiing with data of over 1,00,000 crawled web pages. The crawling was performed and a webgraph was obtained, the data from which gets indexed. Two relevance models namely page rank and HITS are created for ranking the search results after which flat clustering and agglomerative hierarchical clustering is performed. These improve the search results and then query expansion is performed.

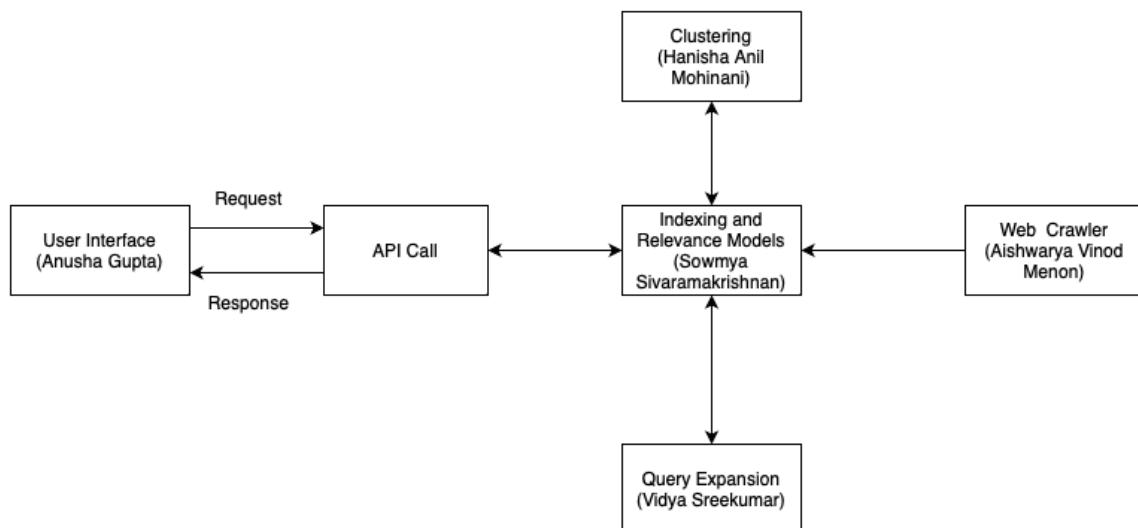
To summarize, we developed an advanced search engine which can be queried to return data about skiing, skiing gear, resorts for skiing, ski festivals, skiing locations etc. using multiple types of relevance models,

clustering techniques and query expansion methodologies, the effectiveness of which are evaluated in the experimental results.

We have divided the development into 5 parts based on different aspects of the project:

- Aishwarya Vinod Menon (Crawling)
- Sowmya Sivaramakrishnan (Indexing and relevance models)
- Anusha Gupta (User Interface)
- Hanisha Anil Mohinani (Clustering)
- Vidya Sreekumar (Query Expansion)

Architecture



Architecture Diagram for Skiing Search Engine

Learnings

- Understanding Information Retrieval Concepts: Developing a search engine requires a solid grasp of information retrieval concepts such as indexing, querying, relevance ranking, and evaluation metrics. Through the process of building a search engine we gained a lot of practical insights into how these concepts are applied in real-world scenarios.
- Hands-on Experience with Search Engine Technologies: Building a search engine involves working with various technologies and tools commonly used in the field, such as web crawlers (eg: Nutch), indexing algorithms like page rank and HITS algorithms, query processing techniques, and relevance ranking algorithms. Experience with these technologies is invaluable for us.
- Overcoming Challenges and Trade-offs: Building a search engine involves making numerous design decisions and trade-offs to balance factors such as search quality, efficiency, scalability, and user experience. We learnt how to overcome the challenges involved in designing and implementing effective search systems, including dealing with large-scale data, optimizing performance, handling query intent ambiguity, and addressing user feedback.

Experience

Working as a project team to develop the search engine was difficult but incredibly gratifying. Collaborating on a project of this nature allowed us to explore the intricacies involved in information retrieval, ranging from developing web indexing algorithms to putting relevance categorization and UI elements into practice. Together, we overcame several technological obstacles, such managing big datasets, optimizing search performance, and debugging intricate code. Nevertheless, overcoming these challenges gave our group a strong feeling of collaboration and improved our problem-solving abilities. Furthermore, it was immensely satisfying to witness the search engine's slow development from a notion to a tool that can find pertinent information on skiing.

Challenges

Developing a search engine with 170,642 indexed web pages and incorporating multiple clustering techniques and relevance models presents several significant challenges.

First and foremost is the sheer scale of the data. Processing and indexing such a large number of web pages requires a solid infrastructure and efficient algorithms to ensure timely search responses. Managing the storage and retrieval of these pages

Aishwarya Vinod Menon (AXV220062), Sowmya Sivaramakrishnan (SXS230043), Anusha Gupta (AXG230026), Hanisha Anil Mohinani (HXM220089), Vidya Sreekumar (VXS22066)

while maintaining retrieval performance can strain computing resources and require careful optimization.

Second, the integration of multiple clustering techniques increases system complexity. Each clustering method may have its own parameters, assumptions, and computational requirements. Coordinating these techniques to effectively organize search results can be difficult, especially when considering different types of ski-related content such as gear, resorts, techniques and news.

Additionally, incorporating different relevance models adds complexity. Relevance models determine how search results are ranked and presented to users based on their query intent. Developing and refining these models to precisely match user preferences and information needs for ski-related queries requires extensive testing and evaluation. In addition, ensuring the accuracy and relevance of search results is a major challenge. The quality, relevance and reliability of ski-related content can vary greatly from site to site. Filtering out irrelevant or low-quality content and promoting authoritative and useful resources requires complex ranking algorithms and constant monitoring and adjustment.

2. Crawling

Web crawling is the process of automatically collecting data from websites. To do this we use web crawlers or spiders which search the internet and get the relevant websites related to the one we are crawling.

First, we provide a list of seed URLs to the web crawlers and then the crawlers visit the hyperlinks mentioned in those websites. Web crawlers continue this process iteratively until we have visited all the seed URLs.

Web crawlers can collect text, images, hyperlinks, videos and metadata from these sites and we can use these to build Search Engine, SEO optimization, Price Monitoring and Market Research.

An open source framework called Apache Nutch is used to create search engines and index the web. This targets a collection of seed URLs that serve as the foundation for research and starts the indexing process. Nutch uses its search functions to effectively get the related web pages by sending HTTP requests to these URLs. During the search process, the framework looks at these pages' content structures and uses a variety of plugins to extract useful data including text, links, and metadata.

A total of 93 seed URLs were gathered by visiting relevant skiing sites and pasted into a text file and passed to Nutch, from which we were able to generate 124659 webpages.

Seed URLs

<https://www.sunvalley.com/>

<https://www.brettonwoods.com/>

<https://en.wikipedia.org/wiki/Skiing>

<https://www.skicanada.org/ready/what-is-skiing/>

<https://www.britannica.com/sports/skiing>

<https://www.skimag.com/>

<https://www.onthesnow.com/>

<https://www.theskibum.com/ski-snow/ski-equipment>

<https://www.sunandski.com/c/snow-ski-equipment>

<https://shop.rockymountainsskiandboard.com/>

https://www.blackdiamondequipment.com/en_US/shop/ski-snowboard/

<https://www.amazon.com/ski-gear/b?ie=UTF8&node=2342470011>

<https://freeskier.com/>

<https://thesnowmag.com/>

<https://www.fall-line.co.uk/>

https://magazines.feedspot.com/ski_magazines/
<https://alltracksacademy.com/blog/the-top-5-skiers-in-the-world/>
<https://snowbrains.com/the-10-greatest-alpine-skiers-of-all-time/>
<https://www.sportsmanskihaus.com/>
<https://www.skitalk.com/forums/>
<https://snowheads.com/ski-forum/>
<https://www.newschoolers.com/forum>
<https://www.tetongravity.com/forums/forumdisplay.php/3-Ski-Snowboard>
<https://www.the-house.com/ski>
<https://justacoloradogal.com/skier-slang-decoded/>
<https://www.nbs.org/>
https://en.wikipedia.org/wiki/Alpine_skiing
<https://tickettoridegroup.com/blog/why-skiers-snowboarders-make-the-best-mates/>
<https://planetski.eu/>
<https://www.yardbarker.com/skiing>
<https://www.snowindustrynews.com/>
<https://www.eurosport.com/alpine-skiing/>
<https://apnews.com/hub/skiing>
<https://www.rei.com/learn/expert-advice/how-to-ski.htmlhttps://www.sunvalley.com/>
<https://www.brettonwoods.com/>
<https://en.wikipedia.org/wiki/Skiing>
<https://www.skicanada.org/ready/what-is-skiing/>
<https://www.britannica.com/sports/skiing>
<https://www.skimag.com/>
<https://www.onthesnow.com/>
<https://www.theskibum.com/ski-snow/ski-equipment>
<https://www.sunandski.com/c/snow-ski-equipment>
<https://shop.rockymountainsskiandboard.com/>
https://www.blackdiamondequipment.com/en_US/shop/ski-snowboard/
<https://www.amazon.com/ski-gear/b?ie=UTF8&node=2342470011>
<https://freeskier.com/>
<https://thesnowmag.com/>
<https://www.fall-line.co.uk/>
https://magazines.feedspot.com/ski_magazines/
<https://alltracksacademy.com/blog/the-top-5-skiers-in-the-world/>
<https://snowbrains.com/the-10-greatest-alpine-skiers-of-all-time/>
<https://www.sportsmanskihaus.com/>
<https://www.skitalk.com/forums/>
<https://snowheads.com/ski-forum/>
<https://www.newschoolers.com/forum>
<https://www.tetongravity.com/forums/forumdisplay.php/3-Ski-Snowboard>
<https://www.the-house.com/ski>
<https://justacoloradogal.com/skier-slang-decoded/>
<https://www.nbs.org/>

https://en.wikipedia.org/wiki/Alpine_skiing

https://tickettoridegroup.com/blog/why-skiers-snowboarders-make-the-best-mates/

https://planetski.eu/

https://www.yardbarker.com/skiing

https://www.snowindustrynews.com/

https://www.eurosport.com/alpine-skiing/

https://apnews.com/hub/skiing

https://www.rei.com/learn/expert-advice/how-to-ski.html

https://www.sundayriver.com/learning-to-ski-and-snowboard

https://thepointsguy.com/guide/right-age-start-kids-skiing/

https://www.googleadservices.com/pagead/aclk?sa=L&ai=DChcSEwjwmtKzwcqFAxXYONQBHX8ZCA4YABAAGgJvYQ&ase=2&gclid=CjwKCAjw5v2wBhBrEiwAXDDoJRukByqSP-Si7g9QM3tiMUTQBGNU0NpWQGj4lpR8vMJakPx70ju6RBoCrpEQAvD_BwE&ohost=www.google.com&cid=CAESVuD2fZ7SDkADMhPr5uRWLnSDbTlu1NOKECIsjKABYAP84eCyeizXRhrcC9yX-bS4CZ0sdLyiSUqecWEymh3diV6d8Mct7F3RULNCItFKL-QQ_kfMluaa&sig=AOD64_29QpNhaonG5owqrNU4h9yspuhwAg&q&nis=4&adurl&ved=2ahUKEwjrxcmzwcqFAxVj6skDHQqhCSgQ0Qx6BAgJEAE

https://unofficialnetworks.com/category/news-3/

https://fasterskier.com/

https://www.skiutah.com/blog/categories/news

https://www.denverpost.com/tag/skiing/

https://www.parkrecord.com/news/final-week-of-ski-season-begins-in-park-city/

https://www.forbes.com/sites/wendysachs/2019/01/14/women-are-transforming-the-ski-industry-on-the-slopes-and-in-the-boardroom/

https://www.ridestore.com/mag/female-skiers-to-follow-on-instagram/

https://www.theskidiva.com/forums/index.php?threads/skiings-gender-gap-widens-with-age.28264/

https://welove2ski.com/how-to-ski/top-olympic-male-skiers/

https://en.wikipedia.org/wiki/List_of_FIS_Alpine_Ski_World_Cup_men%27s_race_winnners

https://www.ranker.com/list/famous-male-alpine-skiers/reference

https://www.ranker.com/list/famous-male-alpine-skiers/reference

https://www.onthesnow.com/

https://travel.usnews.com/rankings/best-ski-destinations/

https://www.forbes.com/uk/advisor/travel-insurance/worlds-best-ski-resorts/

https://www.travelandleisure.com/best-ski-resorts-in-the-world-8406637

https://www.lakeplacid.com/do/outdoors/winterspring/skiing-riding

https://www.denverpost.com/tag/ski-accidents/

https://www.lonestarskiers.com/

https://www.metroplexskiclub.com/

https://www.texasskirangers.org/

https://www.dallaskiclub.org/content.aspx?page_id=22&club_id=592979&module_id=177125

https://en.wikipedia.org/wiki/Cross-country_skiing
<https://www.curated.com/journal/20000/a-guide-to-the-different-kinds-of-skis>
https://buckmans.com/blog/the-complete-guide-to-different-types-of-skis-bpid_445.aspx
<https://www.mountainheaven.co.uk/blog/an-introduction-to-different-types-of-skis/>
<https://www.onthesnow.com/news/which-ski-is-right-for-you/>
<https://www.rei.com/learn/expert-advice/alpine-ski-bindings.html>
<https://www.undercovertourist.com/blog/best-time-ski/>
<https://www.thorindustries.com/stories/skiing-north-america-year-round>

Nutch Crawling Process

The following steps are followed for crawling:

1. Collecting seed URL list and inject it to CrawlDB database.
2. Produce a fetch list from the CrawlDB database and upload this fetch list in a new segment directory.
3. On this new segment directory, fetch the data from the seed URLs.
4. Parse the contents of the seed URLs.
5. Update CrawlDB with the new set of URLs.
6. Repeat steps 2-5 for N number of iterations to get the desired number of web pages crawled. (For our project, we crawled 124,659 webpages).

Components of Nutch Data

CrawlDB contains information about all the URLs, including whether it was fetched or not.

LinkDB contains the links of the known seed URLs which includes the source URL and anchor text of the link.

Segments are a set of URLs fetched together as a unit. It consists of

- crawl_generate - Set of URLs that needs to be fetched.
- crawl_fetch - fetch status of each URL
- Content - raw content got each URL
- parse_text - parsed text of each URL
- parse_data - outlinks and metadata parsed from each URL
- crawl_parse - outlink URLs which are used to update the crawlDB

Below is the Nutch command we used to automate the crawling process to get 124,659 webpages as a result iteratively for 17 iterations.

```
bin/crawl -i -D solr.server.url=http://localhost:8983/solr/nutch -s urls crawl 17
```

Duplicate Results Handling

The CrawlIDB component of the Nutch can easily detect if the URL is already parsed during the update CrawlIDB command execution after each iteration. For the duplicate results which have similar content but different URL is handled by Nutch script which internally runs the “dedup” MapReduce job which marks these results in the CrawlIDB. These marked results are later deleted in the Solr during the indexing process.

Shortest URL, highest score or most recently fetched are some of the heuristics we can use to choose the item which is not marked as duplicate.

Collaboration with Indexing and Relevance Models

The process involved setting up the crawling infrastructure, which included the installation and configuration of Nutch and Solr, as they are interconnected components. The functionality of Nutch's automated script, linked with Solr for indexing, was comprehensively discussed. Input of seed URLs and the corresponding output were reviewed collaboratively. Remote access to the server facilitated joint exploration and understanding of the automation script provided by Nutch, which manages the indexing process through Solr.

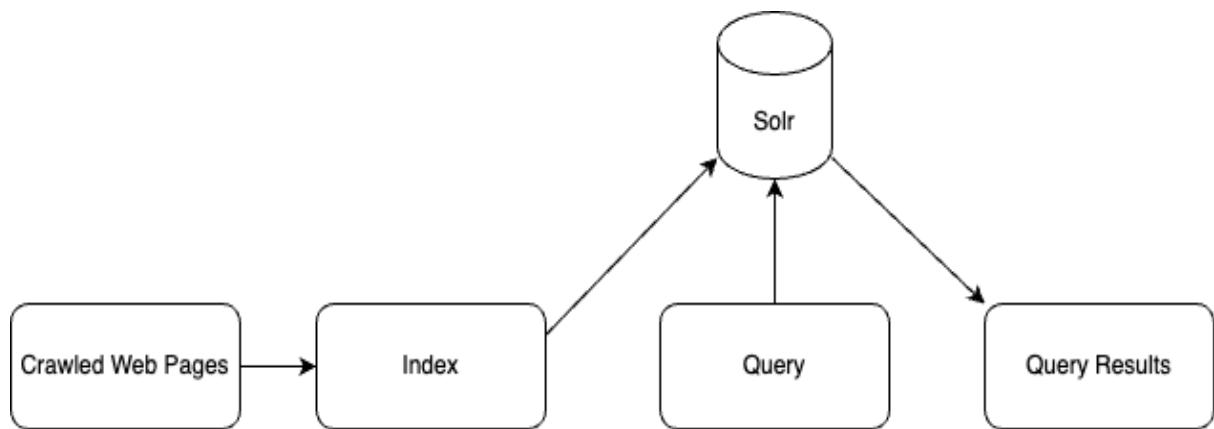
3. Indexing and Relevance

After performing the web crawling with apache nutch, we then proceeded to index the crawled pages using Apache Solr. We configured nutch to internally use solr for indexing and created the index using the following command:

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

In the above command, bin/nutch points to the binary of apache nutch and its an indication that the nutch tool is being utilised to run the command. This is followed by the keyword index which is essentially the command instructing nutch to index the crawled data. We then provide the paths to crawldb- a database storing information about crawled web pages, such as their URLs, metadata, and status and to linkdb- which contains data about the links between pages. The directory where the segments to be indexed are stored is added next in the command. The -filter option indicates that you want to apply any configured URL filters during the indexing process. URL filters can be used to include or exclude certain URLs based on patterns or criteria. -normalize tells Nutch to normalize URLs during the indexing process. URL normalization ensures that URLs are standardized and consistent, which can help improve indexing and reduce duplication of content. Lastly, -deleteGone option instructs Nutch to delete any documents from the index that are marked as "gone" in the crawldb. "Gone" typically refers to web pages that were previously crawled but are no longer accessible or valid.

The diagram below represents flow in the process of assembling an index:



Next we move onto creating web graphs. Creating a web graph in this context helps in understanding and analyzing the structure and organization of content on the websites better. It also facilitates detection of updates made to the pages by drawing comparisons with previously obtained link structure. Web crawler was used to crawl the web pages and extract their urls which are then used by us to create a web

graph. The web pages are each represented as a node of the web graph and the edges connecting these nodes represent the links between them.

The next course of action was to determine the relevance or importance of these web pages for which we explored three different models such as the vector space relevance model, page rank and HITS algorithm.

The first one was the Vector Space Relevance model which is provided by Solr itself. The VSM employed by Solr for information retrieval is pivotal in efficiently assessing the relevance of documents to user queries. Within this model, documents and queries are transformed into vectors within a high-dimensional space, with each dimension corresponding to a unique term in the document collection. Central to this approach is the Term Frequency-Inverse Document Frequency (TF-IDF) weighting scheme, which evaluates the importance of a term within a document relative to its frequency across the entire collection. Documents are represented as vectors, with each dimension reflecting the TF-IDF weight of the corresponding term, while queries are similarly transformed into vector representations based on their constituent terms and TF-IDF weights.

Once documents and queries are represented as vectors, Apache Solr calculates relevance scores using cosine similarity, measuring the cosine of the angle between the query vector and the document vectors. This comparison determines the degree of similarity between the query and each document, with higher cosine similarity scores indicating greater relevance. By ranking documents based on these similarity scores, Solr identifies the most relevant documents to return as search results. Through this Vector Space Model, Apache Solr leverages TF-IDF weighting and cosine similarity to efficiently and effectively retrieve information, facilitating accurate and meaningful search experiences for users across diverse applications.

The next approach was to use the PageRank algorithm for determining the relevance. This algorithm assigns a numerical weight to each element of a hyperlinked set of documents, with the purpose of measuring its relative importance within the set. PageRank interprets links between pages as endorsements or votes for the linked page's importance, with more influential pages receiving higher scores. This iterative algorithm considers the entire link structure of the web, providing a nuanced understanding of relevance beyond individual page content. The built-in nutch webgraph command was employed to create the webgraph from crawled URLs.

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

The PageRank algorithm runs repeatedly until the point of convergence and this was done with the following command:

bin/nutch linkrank -webgraphdb crawl/

With a damping factor of 0.80, we got the page rank algorithm to assign scores to the web pages which was then updated in the crawled database and indexed using Solr. The score update was done with the help of the following command:

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

This is followed by the indexing as mentioned above which is done with the following command:

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

The last algorithm we used was the HITS algorithm for creating a relevance model. This works on determining the relevance on the basis of the authority and hub scores. We wrote our own code to generate these scores based on the inlinks and outlinks we'd obtained earlier for all the web pages. We first identified the most relevant web pages which become our roots, these were identified based on the queries provided by the user. Pages linking to these roots are called hubs while the pages being linked to by the roots are called authorities. We collected a root set of web pages based on page ranks and this was then converted into a base set by adding the inlinks and outlinks. Then we calculated the HIT scores and sorted the max authority scores.

To summarize, different models were used to find the relevance. While each one had its own set of strengths and limitations, we used them specific to the context of our project. The purpose of these relevance models was to give accurate results when the user queries our engine.

Web Graph & Statistics

We created a web graph from the crawled documents. The following command was used to create the web graph from all the crawled web pages:

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

This command processes multiple segments of crawled data together and generates a web graph in different components to a given directory. The three components created are:

- InLink Database: List of URLs and all their inlinks.

- OutLink Database: List of URLs and their outlinks.
- Node Database: List of URLs with meta information including the number of inlinks and outlinks.

Statistics

Total Number of Links: 124659

Total Number of Nodes: 123668

Max number of incoming links: 44

Max number of outgoing links: 2

Web graph information connected to the index

While creating the index in Solr, we have run the LinkRank program of Nutch to perform an iterative link analysis with the web graph. 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.

Highest Hub and Authority Scores

Highest hub score- 0.0027258360565092206

URL with highest hub score- <https://www.skiinfo.no/>

Highest authority score- 0.0027258360565107728

URL with highest authority score- <https://www.skiinfo.no/>

Topic Based Page Ranking

For the purpose of testing topic based page ranking, we chose the topic:

Skiing Jumps

The results from our search engine are as follows:

1. Skiing - Freestyle, Tricks, Jumps | Britannica

<https://www.britannica.com/sports/skiing/Freestyle-skiing>

Score- 0.21385658

2. Ski jumps and balls | SkiTalk | Ski reviews, Ski Selector

<https://www.skitalk.com/threads/ski-jumps-and-balls.33349/>

Score- 0.21375658

3. Jacob_W - Videos - Newschoolers.com

<https://www.newschoolers.com/member/Jacob-W.5247/Videos>

Score- 0.16821946

4. U.S. ski team's 5 essential ski training exercises - OnTheSnow

<https://www.onthesnow.com/news/us-ski-teams-training-exercises/>

Score- 0.15607622

5. SRAM jumps into the ebike powertrain game | SkiTalk | Ski reviews, Ski Selector
<https://www.skitalk.com/threads/sram-jumps-into-the-ebike-powertrain-game.33579/>
Score- 0.15000461

Testing with Queries

The relevance models were tested using 60 queries which we obtained by looking into the most frequently searched queries with respect to skiing on google. Upon obtaining the search results we manually inspected the results by comparing the top results obtained for each query with their page rank or hits scores, depending on the model being tested, and the search results were both relevant to the query and ordered according to their relevance scores.

Collaboration with UI and test relevance models results

We generated several queries to test the various relevance models implemented in this search engine. I have generated 60 queries and tested them on the UI by differentiating the results of PageRank and HITS relevance model results. Upon obtaining the search results we manually inspected the results by comparing the top results obtained for each query with their page rank or hits scores, depending on the model being tested, and the search results were both relevant to the query and ordered according to their relevance scores. The results were judged in some ways such as the HITS relevance model giving the most relevant results in comparison of page rank.

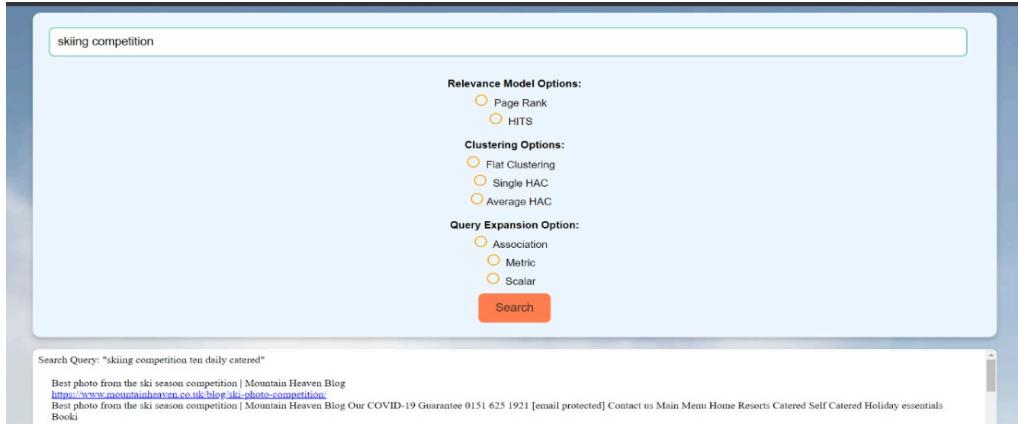
Collaboration with Clustering

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

4. User Interface and Comparison with Google and Bing:

Design of the Interface

I have used HTML, CSS and Javascript to create the front end of the search engine. Additionally, I have used the Flask API to connect it to the backend and then call relevant functions to display the results.



Front end consists of a search bar, ranking algorithm options, submit button, frames to display results of Skiing Search Engine, Google and Bing.

To search a query, the user has to do 3 steps:

- Enter the query in the search bar
- Select one of the following options to rank the result: Page Rank, Hits [Relevance Model]; Flat, Single HAC clustering, Average HAC clustering [Clustering Model Options]; Association, Metric and Scalar [Query Expansion Model options].
- Click on the Submit button.

When the user clicks on submit, a GET API is called with the searched query and selected ranking option as parameters to the flask backend app. The backend app connects to Solr to find relevant results. The results are then sent to the selected ranking model and then ranked accordingly. The ranked results in JSON format are returned and displayed in the front end. For the same query, Google and Bing results are also displayed in their respective frames.

I have also designed a backend app with Python and Flask that does the following:

1. It receives query and rank type parameters from the front end.
2. Preprocess the user query by removing stopwords and punctuations and sends it to Solr API.
3. Receives data from Solr API and calls appropriate ranking model to rank the Solr results.

4. JSONify the results and send it to the front end.

Accessing the Relevance Model

Sowmya was responsible for indexing and relevance models PageRank and HITS. Sowmya provided me PageRank and HITS scores for crawled URLs. In the backend, I have written code to sort the Solr results based on scores provided by Sowmya for PageRank and HITS relevance models.

Queries Tested

I tested 10 common queries related to skiing with Sowmya and checked if sorted results are being returned as per the scores she provided. Additionally, I tested 10 most searched skiing queries as per Google to check if all the ranking models are working correctly and correct data is being sent via API.

Accessing the Clustering Model:

Hanisha was responsible for clustering. She created functions for flat, Single HAC and Average HAC and I integrated it with the Flask application. The functions are called whenever the user chooses clustering model options. Accordingly results are returned and displayed on the front end. We both tested about 15 queries if it was returning relevant results.

Comparison to Google and Bing:

As compared to Google and Bing, our results are not that good. The primary reason for this is that Google and Bing have more computing resources which has allowed them to crawl much more webpages as compared to ours. Plus, they have gathered the user data since decades and improved upon the result. They also use realtime, state of the art algorithms and have high computational power to process the data giving far more accurate results.

Clustering component:

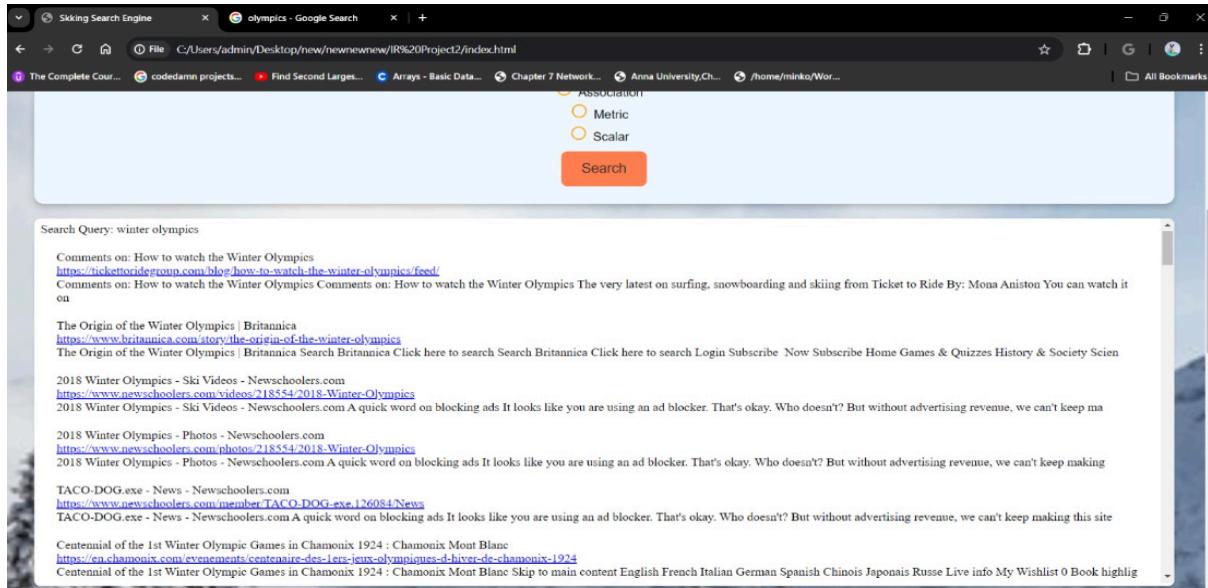
When the user chooses a clustering option, backend code calls a clustering algorithm and returns appropriate results to the front end. The clustering algorithm rearranges the result based on clusters containing top results.

Selecting the queries:

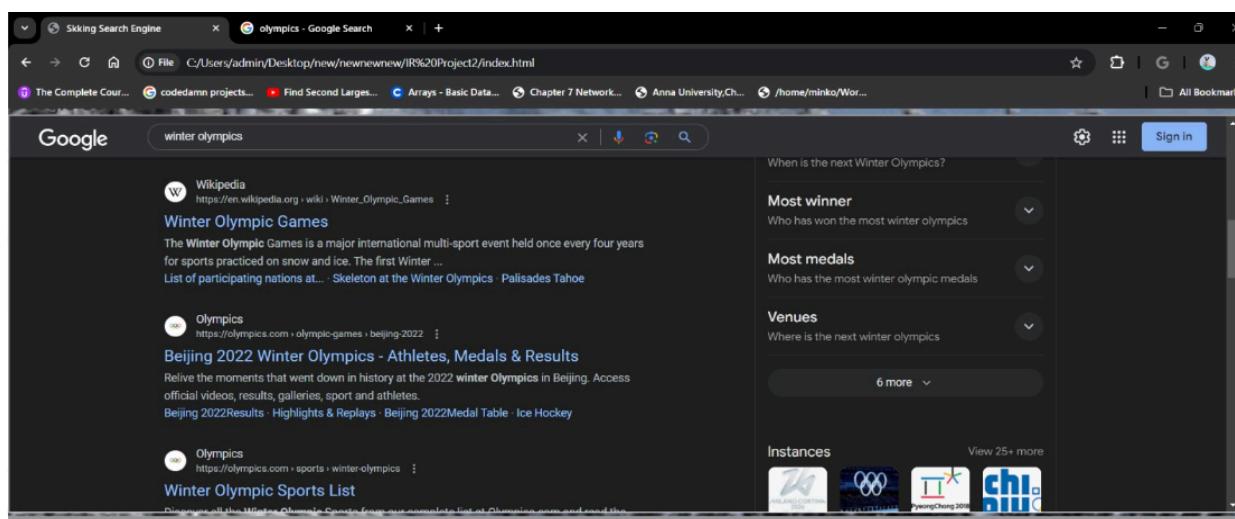
I googled most searched skiing related queries and from those chose the ones which were giving more accurate results in our search engine. We also checked if the results were being improved after clustering and query expansion.

Query Results:

Query 1: Winter Olympics with Page Rank



Skiing Search Engine Result



Google Results

Bing Search results for "olympics".

- Olympics** - https://olympics.com/en/olympic-games/beijing-2022 - Web
- Beijing 2022 Winter Olympics - Athletes, Medals** - whereig.com
- Beijing 2022 Winter Olympics Schedule, Dates, Calendar ...** - whereig.com
- Beijing 2022 Winter Olympics: Stars to watch, event ...** - olympics.com
- Winter Olympic Sports List | Olympics.com** - web

The Winter Olympic Games are a major international multi-sport event held once every four years for sports practiced.

Bing Results

Query 2: Skiing Gear with Flat Clustering

Search Query: skiing gear

- Skiing 101 | BLISTER Outdoor Media & Gear Reviews** - https://blisterreview.com/category/gear-101/skiing-101/
- Skiing 101 | BLISTER Outdoor Media & Gear Reviews** Skip to content Menu Reviews GEAR 101 Skiing 101 Mountain Biking 101 Trail Running 101 Outerwear 101 Boot Fitting 101 Climbing 101
- Backcountry Skiing for Beginners: The Gear - Just a Colorado Gal** - https://justacoloradogirl.com/backcountry-skiing-for-beginners-gear/
- Backcountry Skiing for Beginners: The Gear - Just a Colorado Gal** Home Start Here About Must-Have Gear Backpacking 101 Travel Contact Backcountry Skiing for Beginners: The Gear Share Tweet Pin Missed Ba
- Atomic Skiing Gear | Curated.com** - http://www.curated.com/c/skiing/b/atomic
- Atomic Skiing Gear | Curated.com** DSA_CATEGORY:SKIING DSA_MARKETABLE_CATEGORY:SKIING DSA_DEPARTMENT:WINTER-SPORTS DSA_IS_APPAREL:FALSE Free shipping on orders over \$50 Free personalized recommendations
- Skiing Reaction Gifs - Ski Gabber - Newschoolers.com** - https://www.newschoolers.com/forum/thread/756360/Skiing-Reaction-Gifs
- Skiing Reaction Gifs - Ski Gabber - Newschoolers.com** A quick word on blocking ads It looks like you are using an ad blocker. That's okay. Who doesn't? But without advertising revenue, we can't keep ma
- Inbound Skiing Backpack Contents - Ski Gabber - Newschoolers.com** - https://www.newschoolers.com/forum/thread/873320/Inbound-Skiing-Backpack-Contents
- Inbound Skiing Backpack Contents - Ski Gabber - Newschoolers.com** A quick word on blocking ads It looks like you are using an ad blocker. That's okay. Who doesn't? But without advertising revenue, we c
- Blizzard Skiing Gear | Curated.com** - http://www.curated.com/c/skiing/b/blizzard

Skiing Search Engine Results

Google Search

skiing gear

Filter by r/Skigear

Price Under \$150 \$150 - \$350 Over \$350

In what order do people buy **ski gear**? : r/Skigear - Red... Jan 23, 2020

New skier here - how much should I budget for **gear**? : r/Skigear - Red... Jun 23, 2020

Will I regret getting budget **ski gear** for resort skiing? : r/Skigear - Red... Oct 10, 2020

FAQs About **Gear** and Some Tips for Beginners : r/skiing - Red... Aug 30, 2020

More results from www.reddit.com

Department Men's Women's Kids' Boys' Unisex

Sun & Ski Sports https://www.sunandski.com › Gear & Apparel › Snow ...

Snow Ski Gear - Clothing - Equipment at ... Shop the latest snow ski equipment and apparel at Sun & Ski. Find skis, boots, bindings, and high-quality ski clothing for all your winter adventures!

4.3 ★ store rating (4.1K) · \$136 to \$600 · Free 1–4 day delivery over \$75 · 30-day returns

Stores Backcountry

Google Results

Bing Search

skiing gear

Filter by Kids' Boys' Unisex

Stores Backcountry

REI Co-op https://www.rei.com/h/snowsports +

Skiing & Snowboarding Gear | REI Co-op Find everything you need for your snow adventures at REI, from skis and snowboards to clothing and avalanche safety gear. Explore new arrivals, shop by category, get expert ...

4.3 ★ store rating (4.1K) · \$136 to \$600 · Free 1–4 day delivery over \$75 · 30-day returns

Explore Further

Snowboard Shop - Snowboarding Gear & More | evo evo.com

The Best Winter Sports Gear in 2024 — Ski Gear Brands bestproducts.com

Recommended to you based on what's popular · Feedback

REI Co-op https://www.rei.com/c/ski-clothing +

Ski Clothing | REI Co-op Web Shop for Ski Clothing at REI - Browse our extensive selection of trusted outdoor brands and high-quality recreation gear. Top quality, great selection and expert advice you can trust. ...

50% off 50% off 30% off

Skiing Gear

Map showing locations of Walmart Supercenter, DICK'S Sporting Goods, and FORO SPORTS CLUB in Addison, Mary Kay, Richardson, and Dallas.

Bing Results

Query 3: Skiing Competition with Metric Query Expansion

The screenshot shows a custom search engine interface. At the top, there is a search bar with the query "Skiing competition ten daily catered". Below the search bar, there are two orange circular icons labeled "Metric" and "Scalar". A large orange "Search" button is positioned below these icons. The main content area displays search results for the query. The results include links to various websites such as Mountain Heaven Blog, SkiSurvey Competition, and SkiTalk. The results are presented in a clean, modern style with a light blue background.

Skiing Search Engine Results

The screenshot shows Google search results for the query "skiing competition". The results are displayed on a standard Google search page with a white background. The first result is a link to Wikipedia's "Downhill (ski competition)" page. The second result is a link to the Olympics website. The third result is a link to Salomon's page on ski disciplines. The results are presented in a clean, modern style with a white background.

Google Results

The screenshot shows a web browser window with multiple tabs open. The active tab displays search results for the query "skiing".

Bing Search Results:

- Salomon**
https://www.salomon.com/en-us/what-are-the-disciplines-in-alpine-ski-racing
- What are the disciplines in alpine ski racing?**
Alpine ski racing is organized around six disciplines: Downhill, Super G, Giant Slalom, Slalom, Parallel and Combined. Events are based on speed or/and ...
- Olympics**
https://olympics.com/en/news/olympic-alpine-skling... ▾
- Olympic Alpine skiing at Beijing 2022: Top five things to know**
Web Both Aamodt and Kostelic are the only athletes to have won four gold medals in Alpine skiing at the Winter Olympics. From competition venues to athletes to watch, here are ...
- EXPLORE FURTHER**
- Beijing Olympics: Alpine Courses Still a Mystery | Skiling ...** skilnghistory.org
- Alpine Skiing - Winter Olympics** espn.com
- Recommended to you based on what's popular - Feedback**
- Olympics**
https://olympics.com/en/news/alpine-skling-winter-olympics-sport ▾
- What is alpine skiing? Know all the events and rules - Olympics.com**
Web Overall, alpine skiing has five events - downhill, slalom, giant slalom, super-G and combined. Medals are on offer for both men and women in each event. A mixed team ...
- EXPLORE FURTHER**
- Alpine skiing | History, Events, & Facts | Britannica** britannica.com
- Everything You Need to Know About Alpine Skiing | Ski Judge** skijudge.com
- Recommended to you based on what's popular - Feedback**

See more

Initially, the FIS was only responsible for **Nordic skiing** FIS Nordic World Ski Championships 1925 in Janské Lázně, Czechoslovakia, were given status as the first official World Championships.

In **1924**, at the time of the first Olympic Winter Games, this Commission gave birth to the **Federation Internationale de Ski**

The FIS Freestyle Ski World Cup is an annual freestyle skiing competition arranged by the International Ski Federation since 1980.

Explore more

Bing Results

5. Clustering

Clustering is the process of arranging data items according to how similar they are to one another. The similarity is calculated based on metrics such as Euclidean distance, Cosine similarity, Manhattan distance, etc. The points with the highest similarity score are then grouped together.

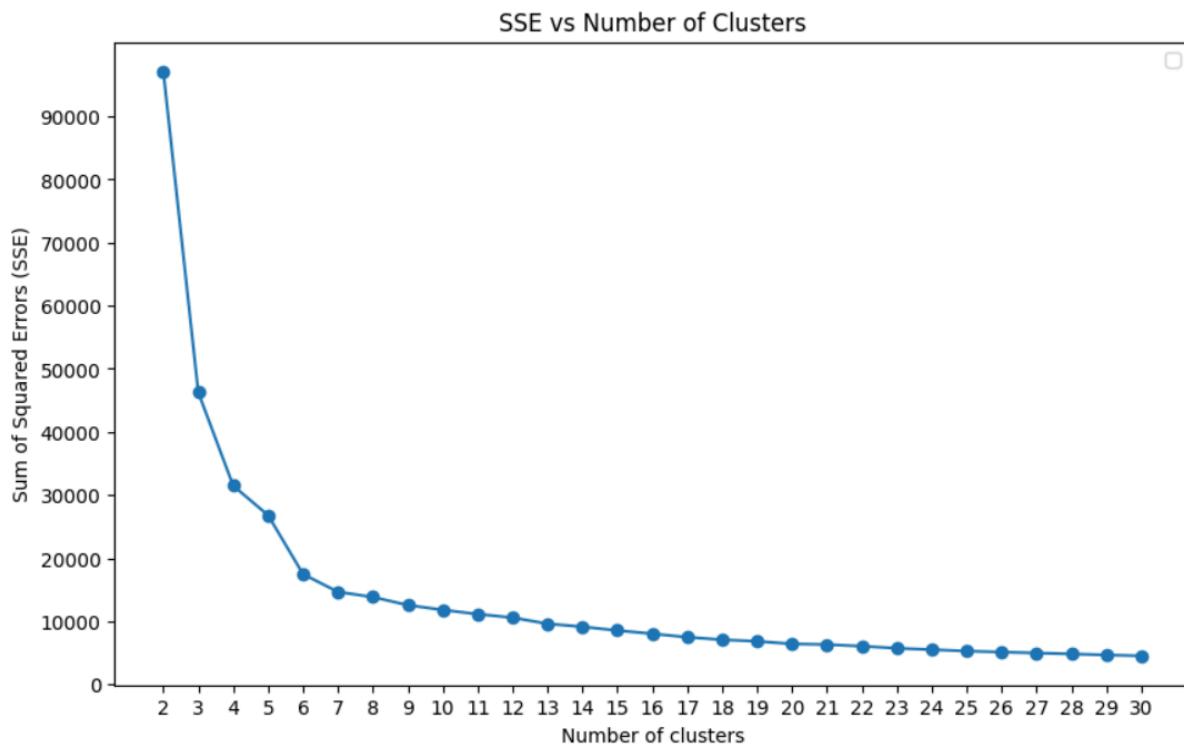
After crawling we received 124,659 webpages. Some pre-processing steps were performed before clustering as given below:

1. The webpages obtained after crawling was in json format
2. The essential features like Title, URL and content were extracted from the file and converted to a pandas dataframe
3. URLs with empty titles or content were removed. Stopwords were also removed
4. The data frame was vectorized using TfidfVectorizer using 10000 features
5. To convert a query to vector notation the tfidf object is saved using pickle

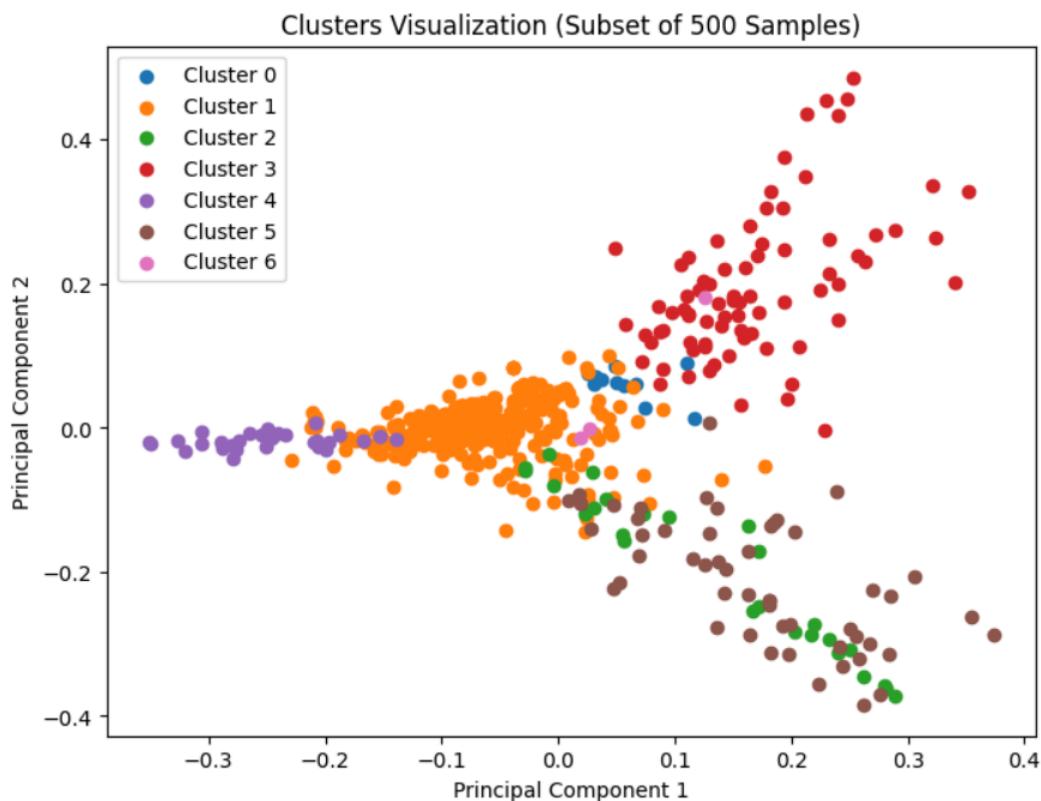
Flat Clustering

Flat clustering was performed using the k-means algorithm. To find the number of clusters required for the k-means algorithm, we used the elbow method. The value of k(cluster number range) was chosen from 2 to 30. The algorithm was applied to each value of k and we calculated the sum of squared distances(SSE) between every data point and its corresponding cluster center. The line chart given below shows the plot of SSE against the number of clusters. Increasing the number clusters decreases the SSE as the distance between points reduces when the number of clusters are increased.

From the line chart we can take k = 7 as the SSE decreases at a slower rate at that point.



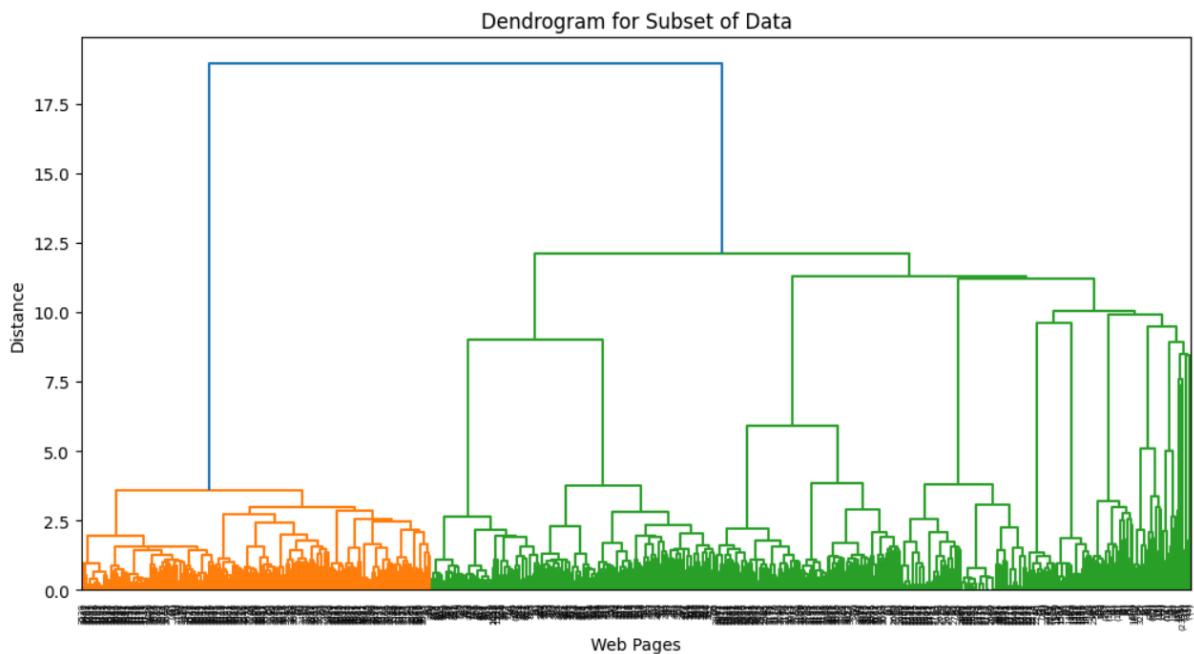
We reduced the dimensionality using Principal Component Analysis(PCA) and plotted 500 samples of the data.



After clustering, a cluster label was assigned to every URL and saved to a file in the format (URL, cluster number). For every cluster, the centroid information was also saved to a file in the format(cluster number, centroid)

Agglomerative Clustering

The same preprocessing steps were performed as that in flat clustering. To determine the number of clusters the hierarchical clustering dendrogram was plotted on the dataset.



Find the highest vertical difference between nodes in the dendrogram, and then pass a horizontal line across the middle. The ideal number of clusters is the number of vertical lines that cross it. In our case we find the number of clusters = 8

Agglomerative clustering was done in batches due to memory constraint. After performing hierarchical clustering with each linkage method (single linkage and average linkage), the centroids of the clusters are calculated. We compare the centroids of each pair of clusters - If the distance between the centroids of two clusters is less than the specified threshold distance of 0.7, the clusters are considered similar enough to be merged. The merging process updates the cluster labels accordingly, ensuring that data points belonging to the merged clusters are assigned the same label.

The above process is repeated until no further clusters can be merged based on the threshold distance and the total number of clusters is 8.

After clustering, a cluster label was assigned to every URL and saved to a file in the format (URL, cluster number). For every cluster, the centroid information was also saved to a file in the format(cluster number, centroid). This was done for both single and average methods in hierarchical clustering.

Integration with front end using Flask API

When the flask app is initialized it reads the 2 files - (URL - Cluster Number) and (Cluster Number - Centroid) and is saved to a dictionary.

Format of results obtained after clustering:

- Input query - All queries are of string type that the user inputs
- Input solr results - Contains the solr results obtained based on the given query
- Output - returns a list of URL objects with its corresponding title, content and other relevant fields

When the user selects a clustering method, the corresponding clustering function with the type of clustering is called. The clustering methods would then find the relevant URLs that match the query that the user had entered and return it to the front-end in JSON format. This would be parsed and results will be displayed in the front-end

Integrate Clustering with Relevance Model

After converting crawled web pages to vector space using Apache SOLR, we integrated related web sites using the hierarchical algometric clustering approach and k-means. The centroid of the cluster was utilized to apply pertinent information. We collect a cluster id, keep track of a list of URLs around the centroid, and utilize that list as our output. To identify the relevant clusters, we utilize the vector space relevance approach for the given query and cluster centroids. We utilize the Page-Rank technique to identify the most important web pages from each cluster after relevant clusters have been extracted. The relevance model is used to identify and elevate the most relevant clusters.

Improving the Search Results using Clustering

- The TfidfVectorizer is used to convert the query entered by the user to a vector form (numerical vector representation).
- The euclidean distance between the query vector and the centroids of all the clusters are calculated.
- The minimum distance gives us the centroid of the cluster that is closest to the query vector.

- Thus all the URLs that belong to that cluster will be returned in JSON format

Results of Flat Clustering

The screenshot shows a search interface with a header containing a 'Scalar' button and a 'Search' button. The search query is 'skiers'. The results list includes various links and text snippets related to skiers, such as gear reviews, photos, and news articles from websites like Skitalk, Newschoolers, and Bob and Terry's.

```

Search Query: skiers
gear for skiers by skiers | SkiTalk | Ski reviews, Ski Selector
https://www.skitalk.com/tag/gear-for-skiers-by-skiers/
gear for skiers by skiers | SkiTalk | Ski reviews, Ski Selector Menu Home Forums New posts Search forums Ski Reviews Articles New articles New comments Search articles eStore Contribute Pugcycle Media
schmuck - Photos - Newschoolers.com
https://www.newschoolers.com/member/schmuck_1702/Photos
schmuck - Photos - Newschoolers.com A quick word on blocking ads It looks like you are using an ad blocker. That's okay. Who doesn't? But without advertising revenue, we can't keep making this site aw
Newschoolers.com
https://www.newschoolers-on-tumblr.com/tagged/videos
Newschoolers.com Ask me anything Submit a post Archive videos The Latest r f t p l Apr 30, 2014 / 25 notes View video page here Video of the Day, April 26th, 2014 // Skiers: Jordan Innes, Matt Wilco
Rental Packages and Pricing — Bob and Terry's Ski and Bike
https://www.bobandterrys.com/rental-pricing
Rental Packages and Pricing — Bob and Terry's Ski and Bike Home Packages & Pricing Tuning Bikes Locations & Hours Book Your Rental Home Packages & Pricing Tuning Bikes Locations & Hours Book Your Rent
Comments on: Why Skiers and Snowboarders make the best mates!
https://tickettoridgegroup.com/blog/why-skiers-snowboarders-make-the-best-mates/feed/
Comments on: Why Skiers and Snowboarders make the best mates! Comments on: Why Skiers and Snowboarders make the best mates! The very latest on surfing, snowboarding and skiing from Ticket to Ride
Comments on: Five of the best resorts for non-skiers
https://www.onthesnow.co.uk/news/best-resorts-for-non-skiers/feed/
Comments on: Five of the best resorts for non-skiers Comments on: Five of the best resorts for non-skiers Enabling the ski travel experience

```

Results of Single Hierarchical Clustering

The screenshot shows a search interface with a header containing a 'Scalar' button and a 'Search' button. The search query is 'skiers'. The results list includes various links and text snippets related to skiers, such as news articles and trail maps from OnTheSnow and Coronet Peak.

```

Search Query: skiers
There's a melting gondola at the top of Aspen Mountain. The artist hopes it will remind skiers about
https://www.cpr.org/2022/04/22/colorado-skiing-aspen-mountain-climate-change-sculpture/
There's a melting gondola at the top of Aspen Mountain. The artist hopes it will remind skiers about the threat of climate change Skip to content News Classical Indie 102.3 KRCC Listen Live Need help?
Five of the best resorts for non-skiers - OnTheSnow
https://www.onthesnow.co.uk/news/best-resorts-for-non-skiers/
Five of the best resorts for non-skiers - OnTheSnow Main Navigation Snow Reports Trip Planning Magazine Search UK EN UK DE FR IT PL SK CZ NO DK NL ES SE Five of the best resorts for non-skiers
Newsoo
Comments on: Five of the best resorts for non-skiers
https://www.onthesnow.co.uk/news/best-resorts-for-non-skiers/feed/
Comments on: Five of the best resorts for non-skiers Comments on: Five of the best resorts for non-skiers Enabling the ski travel experience
Coronet Peak Trail Map | OnTheSnow
https://www.onthesnow.com/new-zealand/coronet-peak/trailmap
Coronet Peak Trail Map | OnTheSnow Snow Reports Trip Planning Magazine US US UK SV SK PL DE ES IT NO DA FR NL CZ Snow Reports Trip Planning Magazine New Zealand , Coronet Peak Related: New Zealand Ota
Share your experience at ski resorts to help other skiers - OnTheSnow
https://www.onthesnow.com/news/share-your-experience-at-ski-resorts-to-help-other-skiers/
Share your experience at ski resorts to help other skiers - OnTheSnow Main Navigation Snow Reports Trip Planning Magazine Search EN EN UK DE FR IT PL SK CZ NO DK NL ES SE Share your experience at ski
Tignes Piste Map | Plan of ski slopes and lifts | OnTheSnow

```

Results of Average Hierarchical Clustering

Search Query: skiers

Five of the best resorts for non-skiers - OnTheSnow
<https://www.onthesnow.co.uk/news/best-resorts-for-non-skiers/>

Five of the best resorts for non-skiers - OnTheSnow Main Navigation Snow Reports Trip Planning Magazine Search UK EN UK DE FR IT PL SK CZ NO DK NL ES SE Five of the best resorts for non-skiers Newsroo

Comments on: Five of the best resorts for non-skiers
<https://www.onthesnow.co.uk/news/best-resorts-for-non-skiers/feed/>

Comments on: Five of the best resorts for non-skiers Comments on: Five of the best resorts for non-skiers Enabling the ski travel experience

Coronet Peak Trail Map | OnTheSnow
<https://www.onthesnow.com/new-zealand/coronet-peak/trailmap>

Coronet Peak Trail Map | OnTheSnow Snow Reports Trip Planning Magazine US US UK SV SK PL DE ES IT NO DA FR NL CZ Snow Reports Trip Planning Magazine New Zealand , Coronet Peak Related: New Zealand Ota

Share your experience at ski resorts to help other skiers - OnTheSnow
<https://www.onthesnow.com/news/share-your-experience-at-ski-resorts-to-help-other-skiers/>

Share your experience at ski resorts to help other skiers - OnTheSnow Main Navigation Snow Reports Trip Planning Magazine Search EN EN UK DE FR IT PL SK CZ NO DK NL ES SE Share your experience at ski

Tignes Piste Map | Plan of ski slopes and lifts | OnTheSnow
<https://www.onthesnow.co.uk/northern-alps/tignes/pistemap>

Tignes Piste Map | Plan of ski slopes and lifts | OnTheSnow Snow Reports Trip Planning Magazine UK US UK SV SK PL DE ES IT NO DA FR NL CZ Snow Reports Trip Planning Magazine Northern Alps , Tignes Rel

Testing Clustering with Queries

60 queries were used to test the clustering methods. The queries were selected based on top google searches according to the topic ‘skiing’. The results were judged manually.

Query list

- Skiing Equipment Checklist
- Skiing Slang Phrases
- Skiing Accident Statistics
- Winter Olympics Skiing
- Skiing Health Benefits
- Skiing Competition Types
- Skier Demographics Analysis
- Skiing Events Calendar
- Skiing Gear Guide
- Skiing History Timeline
- Skiing Skill Levels
- Skiing Locations Worldwide
- Skiing News Updates
- Advanced Skiing Techniques
- Dallas Skiing Trails
- Ski Vacation Packages
- Water Skiing Basics
- Skiing Costs Breakdown
- Freestyle Skiing Tricks

Skiing Medal Winners
Skiing Safety Tips
Skiing Legends Stories
Skiing Competition Rules
Skiing Benefits Overview
Skiing Gear Innovations
Skiing Levels Explained
Skiing Events Schedule
Dallas Ski Resorts
Water Skiing Equipment
Skiing Vacation Deals
Freestyle Skiing Competitions
Skiing Accident Prevention
Olympic Skiers' Profiles
Skiing Gear Maintenance
Skiing Levels Progression
Skiing Events Near Me
Water Skiing Techniques
Skiing Cost Factors
Freestyle Skiing Jumps
Skiing Medalists History
Skiing Gear Reviews
Skiing Levels Comparison
Skiing Events Today
Water Skiing Safety Tips
Skiing Cost Estimate
Freestyle Skiing Tricks List
Skiing Legends Legacy
Skiing Gear Recommendations
Skiing Levels Assessment
Skiing Events Calendar
Water Skiing Locations
Skiing Cost Calculator
Freestyle Skiing Styles
Skiing Legends Stories
Skiing Gear Upgrades
Skiing Levels Evaluation
Skiing Events This Weekend
Water Skiing Gear Guide
Skiing Cost Breakdown
Freestyle Skiing Techniques
Skiing Legends Achievements
Skiing Gear Comparison
Skiing Levels Progress

Skiing Events Nearby
 Water Skiing Basics Guide
 Skiing Cost Analysis
 Freestyle Skiing Competitors
 Skiing Legends Impact
 Skiing Gear Selection
 Skiing Levels Mastery

Query: Skiing Gear

Solr Results

1. <https://www.curated.com/c/skiing/b/atomic>
2. <https://www.curated.com/c/skiing/b/blizzard>
3. <https://www.curated.com/c/skiing/b/fischer>
4. <https://www.curated.com/c/skiing/b/head>
5. <https://www.curated.com/c/skiing/b/k2>

Flat Clustering Results

1. <https://www.skicooper.com/snowsports-safety/>
2. <https://justacoloradogal.com/backcountry-skiing-beginners-gear/>
3. <https://www.newschoolers.com/forum/thread/938284/Most-annoying-term-in-skiing>
4. <https://www.theskidiva.com/>
5. <https://www.newschoolers.com/forum/thread/756360/Skiing-Reaction-Gifs>

Single Hierarchical Clustering Results

1. <https://www.onthesnow.co.uk/news/ski-packing-list/>
2. <https://www.onthesnow.co.uk/news/the-best-sustainable-ski-gear/>
3. <https://www.onthesnow.co.uk/news/how-to-find-the-best-skiing-gloves/>
4. <https://www.onthesnow.com/news/gear/page/2/>
5. <https://www.tetongravity.com/community/login>

Average Hierarchical Clustering Results

1. <https://www.onthesnow.co.uk/news/ski-packing-list/>
2. <https://www.onthesnow.co.uk/news/the-best-sustainable-ski-gear/>
3. <https://www.onthesnow.co.uk/news/how-to-find-the-best-skiing-gloves/>
4. <https://www.onthesnow.com/news/gear/page/2/>
5. <https://www.tetongravity.com/community/login>

Query: Skiing resort

Solr Results

1. <https://www.smuggs.com/pages/welcome/>
2. <https://www.smuggs.com/pages/winter/>
3. <https://www.smuggs.com/pages/winter/activities/cat-trax.php>
4. <https://www.smuggs.com/pages/winter/activities/iceSkating.php>
5. <https://www.smuggs.com/pages/winter/activities/tubing.php>

Flat Clustering Results

1. <https://www.smuggs.com/pages/welcome/>
2. <https://www.smuggs.com/pages/winter/activities/iceSkating.php>
3. <https://www.smuggs.com/pages/winter/>
4. <https://www.smuggs.com/pages/winter/amenities/massage.php>
5. <https://www.smuggs.com/pages/winter/lodging/resortMap.php>

Single Hierarchical Clustering Results

1. <https://www.onthesnow.co.uk/upper-austria/hinterstoder/skireport>
2. <https://www.onthesnow.co.uk/tyrol/imst/skireport>
3. <https://www.onthesnow.com/australia/dinner-plain/ski-resort>
4. <https://www.onthesnow.com/australia/dinner-plain/skireport>
5. <https://www.onthesnow.co.uk/aosta-valley/chamois/skireport>

Average Hierarchical Clustering Results

1. <https://www.onthesnow.co.uk/aosta-valley/chamois/skireport>
2. <https://www.onthesnow.co.uk/carinthia/gerlitzen/skireport>
3. <https://www.onthesnow.com/michigan/crystal-mountain/reviews>
4. <https://www.onthesnow.co.uk/upper-austria/hinterstoder/skireport>
5. <https://www.onthesnow.co.uk/tyrol/imst/skireport>

Query: Skiing Safety

Solr Results

1. <https://kinderlift.com/>
2. <https://www.skitalk.com/tags/avalanche-safety/>
3. <https://www.mountainheaven.co.uk/blog/the-skiers-responsibility-code/>
4. <https://www.skitalk.com/threads/daughter-is-ready-to-move-off-of-bunny-slopes-but-too-small-for-lifts.33020/>

Flat Clustering Results

1. <https://kinderlift.com/childrens-ski-vests/feed/>
2. <https://cochranskiarea.com/safety/>

3. <https://www.tahoedonner.com/amenities/>
4. <https://www.curated.com/c/winter-sports-backcountry-safety-equipment>
5. <https://www.skicooper.com/snowsports-safety/>

Single Hierarchical Clustering Results

1. <https://www.onthesnow.co.uk/news/family-freeriding/>
2. <https://www.perisher.com.au/reports-cams/cams>
3. <https://www.skitalk.com/forums/general-skiing.5/index.rss>
4. <https://www.perisher.com.au/>
5. <https://www.mtpeter.com/>

Average Hierarchical Clustering Results

1. <https://www.skitalk.com/forums/general-skiing.5/index.rss>
2. <https://www.sugarbush.com/>
3. <https://www.onthesnow.co.uk/carinthia/falkert/skireport>
4. <https://kinderlift.com/>
5. <https://www.palisadestahoe.com/footer/safety>

6. Query Expansion and Relevance FeedBack

Query Expansion is a technique used to enhance the search results using various methods like:

1. Rocchio Algorithm (Relevance Feedback method)
2. Pseudo Relevance Feedback method:
 - a. Association Cluster based Query Expansion
 - b. Metric Cluster based Query Expansion
 - c. Scalar Cluster based Query expansion

1. Rocchio Algorithm (Relevance Feedback method)

It is used to improve the results of the search engine by refining original query by using a set of relevant documents.

$$q' = \alpha \cdot q + \beta \cdot \frac{1}{|Dr|} \sum_{dj \in Dr} dj - \gamma \frac{1}{|Dnr|} \sum_{dj \in Dnr} dj$$

Where,

q – original query

q' – modified query

Dr – set of relevant vectors

Dnr – set of non – relevant/irrelevant vectors

α, β, γ – weights

The values of $\alpha=1$, $\beta=0.85$, $\gamma=0.1$ were set to separate relevant documents from the irrelevant documents.

It was tested against 20 queries. The queries were chosen based on how closely they are related to the topic(skiing) . While a few queries were directly related to skiing, the other words were not so closely associated with skiing.

The top few relevant and irrelevant pages for different queries are given below

Query	Relevant Pages	Irrelevant Pages
skiing	<ul style="list-style-type: none"> • https://thesnowmag.com/category/heli-skiing-articles/heli-alaska/feed/ • https://thesnowmag.com/category/heli-skiing-articles/heli-canada/feed/ • https://www.onthesnow.co.uk/news/year-round-skiing/ • http://www.chamonixskiguide.com/ • https://thesnowmag.com/category/heli-skiing-articles/feed/ 	<ul style="list-style-type: none"> • https://www.fis-ski.com/alpine-skiing • https://thesnowmag.com/category/heli-skiing-articles/heli-new-zealand/ • https://thesnowmag.com/category/custom-skis/ • https://thesnowmag.com/category/park-city/ • https://thesnowmag.com/category/heli-skiing-articles/heli-united-states/feed/
snowboarding	<ul style="list-style-type: none"> • http://www.mechanicsofsport.com/snowboarding/equipment.html • http://www.mechanicsofsport.com/advertising.html • https://www.onthesnow.co.uk/news/best-snowboarding-boots-how-to-find-the-right-pair/ • https://www.onthesnow.co.uk/news/best-snowboarding-rents/feed/ • https://www.onthesnow.com/news/travel-guide-skiing-and-snowboarding-in-canada/feed/ • https://www.curated.com/journal/1178001/ski-size-chart-how-to-size-skis 	<ul style="list-style-type: none"> • https://www.onthesnow.co.uk/news/featured-gear/feed/ • https://www.curated.com/c/ski-binding-s • https://www.curated.com/c/ski-boot-accessories • https://www.curated.com/c/ski-boots • https://www.curated.com/c/ski-poles • https://www.curated.com/c/ski-skins • https://www.curated.com/c/ski-snowboard-helmets • https://www.curated.com/c/ski-snowboard-sunglasses • https://www.curated.com/c/ski_snowboard_bags

	<ul style="list-style-type: none"> https://www.onthesnow.co.uk/news/best-snowboarding-boots-how-to-find-the-right-pair/feed/ https://www.curated.com/c/snowboarding/b/abroarding/b/ones 	https://www.curated.com/c/skiing/b/atomic
mountains	<ul style="list-style-type: none"> https://www.onthesnow.co.uk/villach-skiing-mountains/projected-openings https://www.onthesnow.co.uk/villach-skiing-mountains/open-resorts https://www.britannica.com/browse/Mountains-Volcanoes https://www.skiinfo.pl/villach-skiing-mountains/warunki-narciarskie https://www.skiinfo.dk/villach-skiing-mountains/snerapport https://www.onthesnow.co.uk/villach-skiing-mountains/webcams https://www.onthesnow.co.uk/villach-skiing-mountains/ski-resorts https://www.onthesnow.co.uk/rhodope-mountains/skireport https://www.onthesnow.co.uk/villach-skiing-mountains/projected-closing 	<ul style="list-style-type: none"> https://www.onthesnow.com/rocky-mountains/lift-tickets https://www.aspensnowmass.com/four-mountains/aspen-highlands/lift-status https://www.onthesnow.com/news/snow-science-how-mountains-make-snow/feed/ https://www.onthesnow.co.uk/villach-skiing-mountains/lodging https://www.onthesnow.com/news/how-mountains-of-california-make-snow/feed/ https://www.onthesnow.com/news/how-the-mountains-of-colorado-make-snow/feed/ https://www.aspensnowmass.com/four-mountains/aspen-highlands https://www.onthesnow.com/news/how-mountains-of-california-make-snow/

equipement	<ul style="list-style-type: none"> • https://thesnowmag.com/category/heli-skiing-articles/heli-alaska/feed/ • https://thesnowmag.com/category/heli-skiing-articles/heli-canada/feed/ • https://www.onthesnow.co.uk/news/year-round-skiing/ • http://www.chamonixskiguide.com/ • https://thesnowmag.com/category/heli-skiing-articles/feed/ • https://www.onthesnow.co.uk/news/featured-weather/ 	<ul style="list-style-type: none"> • https://www.fis-ski.com/alpine-skiing • https://thesnowmag.com/category/heli-skiing-articles/heli-new-zealand/ • https://thesnowmag.com/category/custom-skis/ • https://thesnowmag.com/category/park-city/ • https://thesnowmag.com/category/heli-skiing-articles/heli-united-states/feed/ • https://www.skitalk.com/tags/mogul-skiing/
slopes	<ul style="list-style-type: none"> • https://www.onthesnow.com/indiana/lodging • https://www.onthesnow.co.uk/news/best-spring-skiing-best-to-north-facing-ski-slopes/ • https://www.meribell.net/en/practical-information/slopes/ • https://www.onthesnow.com/new-york;brantling-ski-slopes/webcams • https://www.onthesnow.com/new-york;brantling-ski-slopes/trailmap • https://www.onthesnow.com/indiana/perfect-north-slopes/skireport • https://www.onthesnow.com/indiana/perfect-north-slopes/lodging 	<ul style="list-style-type: none"> • https://www.onthesnow.co.uk/news/skiing-in-the-dolomites-alpine-charm-and-jagged-peaks/ • https://www.onthesnow.co.uk/salzburg/zell-am-see-schmittenhoehe/ski-report-reviews • https://www.onthesnow.co.uk/bernese-oberland/saanenmoeser/ski-report-reviews • https://www.onthesnow.co.uk/news/guide-to-ski-insurance • https://www.onthesnow.co.uk/nearby-resorts • https://www.onthesnow.co.uk/news/indoor-skiing-on-revolving-slopes/feed/

slopes	<ul style="list-style-type: none"> • https://www.onthesnow.com/indiana/perfect-north-slopes/historical-snowfall • https://www.onthesnow.com/indiana/perfect-north-slopes/lift-tickets • https://www.ski-saintgervais.com/en/notre-offre-de-forfait-saint-gervais-evasion • https://www.onthesnow.com/indiana/perfect-north-slopes/ski-report-reviews • https://www.skiresort.info/snow-reports/sorted/open-slopes/ 	<ul style="list-style-type: none"> • https://www.onthesnow.co.uk/news/uk-dry-ski-slopes/feed/ • https://www.onthesnow.com/new-york;brantling-ski-slopes/lodging • https://www.onthesnow.co.uk/news/best-spring-skiing-hed-to-north-facing-ski-slopes/feed/ • https://www.onthesnow.co.uk/news/europe-snow-report/feed/
terrain	<ul style="list-style-type: none"> • https://www.aspensnowmass.com/discovers/experiences/guides/guide-to-lessons • https://www.onthesnow.com/news/visitors-choice-award-for-best-terrain-park-keystone/ • https://www.onthesnow.com/news/jackson-hole-rated-best-all-mountain-terrain-resort/ • https://www.onthesnow.com/news/visitors-choice-award-for-best-all-mountain-terrain-telluride/ • https://www.onthesnow.com/alaska/eaglecrest-ski-area/reviews • https://www.onthesnow.com/new-jersey 	<ul style="list-style-type: none"> • https://www.onthesnow.com/new-york/bristol-mountain/ski-resort • https://www.onthesnow.com/oregon/mt-bachelor/trailmap • https://www.onthesnow.com/arizona/arizona-snowbowl/trailmap • https://www.onthesnow.com/colorado/durango-mountain-resort/trailmap • https://www.onthesnow.com/new-york/holiday-valley/trailmap

	<p>y/mountain-creek-resort/trailmap</p> <ul style="list-style-type: none"> ● https://www.onthesnow.com/alberta/sunshine-village/trailmap 	
resorts	<ul style="list-style-type: none"> ● https://www.onthesnow.co.uk/lillehammer-ski-resorts/open-resorts ● https://www.onthesnow.co.uk/lillehammer-ski-resorts/ski-pass ● https://www.onthesnow.co.uk/lillehammer-ski-resorts/projected-closing ● https://www.skiresort.info/snow-reports/filter/open-ski-resorts/ ● https://www.onthesnow.co.uk/andorra/open-resorts ● https://www.onthesnow.co.uk/slovenia/open-resorts ● https://www.onthesnow.com/australia/open-resorts 	<ul style="list-style-type: none"> ● https://www.onthesnow.co.uk/nidwalden/open-resorts ● https://www.onthesnow.co.uk/oberstaufen/open-resorts ● https://www.onthesnow.co.uk/romania/open-resorts ● https://www.onthesnow.co.uk/sauerland/open-resorts ● https://www.onthesnow.co.uk/schwarzwald/open-resorts ● https://www.onthesnow.co.uk/ticino/open-resorts ● https://www.onthesnow.co.uk/toggenburg/open-resorts
gear	<ul style="list-style-type: none"> ● https://justacolorado.com/category/gear/page/2/ ● https://www.newschoolers.com/forum/thread/931417/My-ski-pant-company ● https://www.newschoolers.com/forum/thread/932497/Gear-to-stay-away-from ● https://www.newschoolers.com/forum/thread/938051/Park-Skis 	<ul style="list-style-type: none"> ● https://www.mountainheaven.co.uk/blog/cheap-ski-gear/ ● https://www.newschoolers.com/forum/thread/937016/Misfit-skis-drama ● https://www.newschoolers.com/forum/thread/714366/Atonic-2012-2013-Thread ● Show-me-your-babies">https://www.newschoolers.com/forum/thread/933300>Show-me-your-babies

	<ul style="list-style-type: none"> https://www.newschoolers.com/forum/thread/481906/The-answer-to--what-binding- 	<ul style="list-style-type: none"> https://www.newschoolers.com/forum/thread/935371/Baggy-snowpants https://www.newschoolers.com/forum/thread/926870/Vs-came-in-today https://www.newschoolers.com/forum/thread/927598/Post-a-Pic-of-Your-Camera-Gear-Equipment
avalanche	<ul style="list-style-type: none"> https://www.skitalk.com/tags/avalanche-safety/ http://www.mtavalanche.com/sitemap http://www.mtavalanche.com/avalanche-incidents https://avalancherecord.ca/ https://www.mtavalanche.com/resources https://www.mtavalanche.com/node/31203 http://www.mtavalanche.com/accidents https://www.mtavalanche.com/accidents https://www.mtavalanche.com/glossary 	Nan
freestyle	<ul style="list-style-type: none"> https://newschoolers-on.tumblr.com/tagged/FREESTYLE%20GIF https://newschoolers-on.tumblr.com/tagged/freestyle%20gif 	<ul style="list-style-type: none"> https://newschoolers-on.tumblr.com/tagged/nsvideos https://newschoolers-on.tumblr.com/tagged/photo https://newschoolers-on.tumblr.com/tagged/

	<ul style="list-style-type: none"> • https://newschooler-s-on.tumblr.com/tagged/FREESTYLE • https://newschooler-s-on.tumblr.com/tagged/VOTD 	gged/style%20gif
Snow	<ul style="list-style-type: none"> • https://www.j2ski.com/snow_forecast/France/Les_Menuires_snow.html • https://www.j2ski.com/snow_forecast/France/The_Three_Valleys_snow.html • https://www.j2ski.com/snow_forecast/France/Val_d_Iserine_snow.html • https://www.j2ski.com/snow_forecast/France/Courchevel_snow.html • https://www.j2ski.com/snow_forecast/France/Les_Deux_Alpes_snow.html • https://au.j2ski.com/snow_forecast/France/Saint_Martin_de_Belleville_snow.html 	<ul style="list-style-type: none"> • https://www.j2ski.com/snow_forecast/France/Brides_les_Bains_snow_forecast_48.html • https://www.onthesnow.co.uk/bernese-oberland/habkern/weather • https://thesnowmag.com/category/winter-ski-fashion/luxury-ski-wear/feed/ • https://www.onthesnow.com/wyoming/snow-king-resort/ski-report • https://www.onthesnow.com/vermont/mount-snow/lodging • https://www.j2ski.com/
Skiing techniques	<ul style="list-style-type: none"> • https://www.onthesnow.co.uk/carinthia/gerlitzen/ski-resort • https://www.onthesnow.co.uk/carinthia/gerlitzen/ski-report • http://www.mechanicsofsport.com/skiing/how_to_ski.html • https://www.onthesnow.co.uk/news/ski-instructor-courses-train-qualify-and-work-on-the-slopes/ • <a href="https://www.lakepla 	Nan

	<p>cid.com/story/2023/10/24/top-5-reasons-to-take-your-family-skiing-at-whiteface-mountain</p> <ul style="list-style-type: none"> • https://www.onthesnow.com/news/experience-the-arlberg-austrias-largest-ski-resort/ • https://www.onthesnow.com/news/reasons-to-ski-grand-targhee/ 	
Snowboard tricks	<ul style="list-style-type: none"> • https://www.newscrollers.com/forum/thread/874399/Tame-dog-with-stiff-skis • https://www.onthesnow.co.uk/news/in-door-skiing-the-pros-and-cons/ 	<ul style="list-style-type: none"> • https://www.onthesnow.com/wyoming/snow-king-resort/ski-report • https://www.onthesnow.com/vermont/mount-snow/lodging
Mountain resorts	<ul style="list-style-type: none"> • https://www.onthesnow.com/news/a-truly-all-inclusive-ski-vacation-with-club-med/ • https://www.onthesnow.com/new-york/labrador-mt/lift-tickets • https://www.onthesnow.com/new-york/song-mountain/lift-tickets • https://www.morzin.esourcemagazine.com/_the_source_files_/ • https://www.onthesnow.com/australia/dinner-plain/lift-tickets • https://www.onthesnow.com/british-columbia/powder-king/lift-tickets 	<ul style="list-style-type: none"> • https://www.onthesnow.com/alaska/eaglecrest-ski-area/lift-tickets • https://www.onthesnow.com/alaska/hilltop-ski-area/lift-tickets • https://www.onthesnow.com/alberta/ski-banff-norquay/lift-tickets • https://www.onthesnow.com/alberta/winsport/lift-tickets • https://www.onthesnow.com/argentina/caviahue/lift-tickets • https://www.onthesnow.com/argentina/chapelco/lift-tickets • https://www.onthesnow.com/argentina

		<ul style="list-style-type: none"> /las-lenas/lift-tickets https://www.onthesnow.com/arizona/arizona-snowbowl/lift-tickets
Powder snow	<ul style="list-style-type: none"> https://www.skiguidejapan.com/EN/backcountry/index.html https://www.skiguidejapan.com/EN/area/geto/geto_resort/index.html https://www.skiguidejapan.com/EN/area/appi_hachimantai/hachimantai_shimokura/index.html https://www.skiguidejapan.com/EN/accommodation/index.html https://www.newschoolers.com/member/iskiatstoke.59236/Videos https://www.skiguidejapan.com/EN/index.html https://es.skiinfo.com/colorado/steamboat/webcams 	<ul style="list-style-type: none"> https://www.onthesnow.co.uk/tyrol/venet/skireport https://www.onthesnow.co.uk/veneto/cortina-dampezzo/ski-report-reviews https://www.onthesnow.co.uk/vorarlberg/gargellen/ski-resort https://www.onthesnow.co.uk/vorarlberg/gargellen/ski-report https://www.onthesnow.co.uk/vorarlberg/golm/skireport https://www.onthesnow.co.uk/styria/gaisbergalm/ski-resort https://www.onthesnow.co.uk/styria/loser-sandling-altaussee/ski-resort
Slope conditions	<ul style="list-style-type: none"> https://www.onthesnow.co.uk/brescia/open-resorts https://www.onthesnow.co.uk/highland/open-resorts https://www.onthesnow.co.uk/monterosa/open-resorts https://www.onthesnow.co.uk/rhoen/open-resorts https://www.onthesnow.co.uk/abruzzo/ 	<ul style="list-style-type: none"> https://www.onthesnow.co.uk/beskid-slaski/open-resorts https://www.onthesnow.co.uk/czarna-gora/open-resorts https://www.onthesnow.co.uk/devoluy/open-resorts https://www.onthesnow.co.uk/eastern-slovakia/open-resorts https://www.onthesnow.co.uk/val-d-isere/val-d-isere

	<ul style="list-style-type: none"> • open-resorts • https://www.onthesnow.co.uk/apuseni/open-resorts • https://www.onthesnow.co.uk/aragon/open-resorts • https://www.onthesnow.co.uk/belgium/open-resorts 	now.co.uk/hautes-pyrenees/open-resorts
Terrain parks	<ul style="list-style-type: none"> • https://www.ullr.ski/ • https://www.ullr.ski/pro/lifts/index.html • https://www.wachusett.com/events-park/terrain-parks/ • https://www.onthesnow.com/wisconsin/little-switzerland/reviews 	<ul style="list-style-type: none"> • https://www.newshoolers.com/forum/thread/925284/TERM-END-WATCH-21--22 • https://www.newshoolers.com/forum/thread/937430/Rails-are-shrinking--- • https://www.newshoolers.com/forum/thread/938507/Is-Eileen-Gu-an-industry-plant-
Skiing gear	<ul style="list-style-type: none"> • https://www.curated.com/c/skiing/b/atomic • https://www.curated.com/c/skiing/b/blizzard • https://www.curated.com/c/skiing/b/fischer 	<ul style="list-style-type: none"> • https://www.onthesnow.co.uk/devoluy/open-resorts
Avalanche safety	<ul style="list-style-type: none"> • https://www.skitalk.com/tags/avalanche-safety/ • https://www.mountainheaven.co.uk/blog/avalanche-safety-what-you-need-to-know-to-stay-safe/feed/ • https://avalancheresearch.ca/ • https://mountaineer.com/activities/climbing/crampons-and-spares-parts/ • https://mountaineer.com/activities/climbing/crash-pads/ • https://mountaineer.com/activities/climbing/dogbones-sewn-slings-and-lanyards/ 	

	<ul style="list-style-type: none"> .com/sports/backcountry-ski/avalanche-safety/ https://caicssignup.ca/avalancheresearch.ca/en/ https://www.mountainheaven.co.uk/blog/ski-schools-in-france-trained-for-avanche-safety/ 	<ul style="list-style-type: none"> https://mountaineer.com/activities/climbing/harnesses/
Skiing holidays	<ul style="list-style-type: none"> https://www.skisolutions.com/ https://www.whitetracks.co.uk/Helicopter_Transfers_Flims_Laax_Switzerland.htm https://www.onthesnow.co.uk/styria/ramsau-am-dachstein/ski-resort https://www.onthesnow.co.uk/styria/ramsau-am-dachstein/skireport https://hu.ski/about/ 	<ul style="list-style-type: none"> https://www.tyrol.com/things-to-do/sports/cross-country-skiing https://www.mountainheaven.co.uk/blog/ski-jackets-the-best-names-to-look-out-for/ https://www.mountainheaven.co.uk/quick-links/ski-jobs

Testing queries and getting their modified Rocchio query

Query	Rocchio Modified Query
skiing	Skiing resort
snowboarding	Snowboarding gear
mountains	Mountain vermont skiing resort
men/women	Men ski apparel jacket
slopes	Slopes skiing trails snow quality
terrain	terrain ski mountain snow sports
resorts	resorts ski lodges mountain amenities

gear	Skiing gear cost
avalanche	avalanche safety precautions skiing
freestyle	Freestyle skiing techniques
snow	Snow conditions skiing
Skiing techniques	Skiing techniques advanced slope
Snowboard tricks	Snowboard tricks olympics
Mountain resorts	Mountain resorts budget secluded
Powder snow	Powder snow conditions report tips
Slope conditions	Slope conditions alps update
Terrain parks	Terrain parks safety location
Skiing gear	Skiing gear cost sale
Avalanche safety	Avalanche safety protocols tips
Skiing holidays	Skiing holidays resorts accommodation

Observation:

1. No Irrelevant Web Pages for Some Queries

For certain queries like “avalanche”, “skiing techniques” showed no irrelevant web pages.

2. Relevance of Irrelevant Web Pages

For some queries, like “freestyle” the web pages identified as irrelevant seemed far more or equally relevant as the pages identified as relevant.

3. We were unable to arrive at a perfect assignment of weights for our queries.

Considering these issues, we decided not to incorporate the Roccio algorithm into the query expansion.

Instead, we identify relevant documents based on the initial query by dividing the top 50 Solr results into two sets: the relevant set of 30 documents and the irrelevant set containing 20 documents. Then we calculate the cosine similarity between the terms in the query and the terms in the relevant documents. Terms from the relevant documents that are similar to the terms in the query are selected and added to the query to expand it.

2. Pseudo Relevance based Query Expansion:

a. Association Cluster based Query Expansion:

It is an expansion technique that is used to enhance the search query by identifying the most frequently co-occurring stems in the relevant documents and works on the underlying principle that these co-occurring terms are associated.

b. Metric Cluster based Query Expansion:

The metric cluster based query expansion is based on correlation between the terms or stems in this case. It identifies clusters based on the proximity of the terms with each other.

c. Scalar Cluster based Query Expansion

Scalar cluster based Query expansion selects the terms based on the co-occurrence patterns and groups terms together that have similar contextual neighborhoods.

The following is the list of 60 queries to test using pseudo relevance feedback

Query	Expanded Query	Relevant pages fetched	Method used
Skiing	Skiing level winter magazine	45	Association
Skiing equipment checklist	Skiing equipment checklist level admin commonly	30	Metric
Skiing slang phrases	Skiing slang phrases part pass year	25	Scalar
Skiing Accident Statistics	Skiing Accident Statistics earth pass souls	13	Association
Winter Olympics Skiing	Winter Olympics Skiing tricia earth turn	22	Metric
Skiing Competition Types	Skiing Competition Types turn run earth	17	Association
Skiing Gear	Skiing Gear like	12	Association

	wood tricia		
Skiing History Timeline	Skiing History Timeline turn organisation seeing	8	Metric
Skiing Locations	Skiing Locations run reflects directions	25	Association
Skiing News	Skiing News resort inches humble	16	Metric
Advanced Skiing Techniques	Advanced Skiing Techniques run levels garbanzo	32	Scalar
Skiing Trails	Skiing Trails heritage direction pass	22	Association
Ski Vacation Packages	Ski Vacation Packages run directions unplugged	11	Metric
Water Skiing	Water Skiing turn run george	31	Scalar
Skiing Costs	Skiing Costs profit turn bonjour	24	Association
Freestyle Skiing	Freestyle Skiing potd candles epic	30	Metric
Skiing Medal Winners	Skiing Medal Winners imperial ranks meribels	19	Scalar
Skiing Safety Tips	Skiing Safety Tips turn heritage directions	9	Association
Skiing Legends Stories	Skiing Legends Stories turn interviews local	21	Metric
Skiing Benefits Overview	Skiing Benefits Overview turn directions policy	14	Association
Skiing Levels Explained	Skiing Levels Explained	28	Scalar
Skiing Events Schedule	Skiing Events Schedule run pass local	32	Association
Ski Resorts	Ski Resorts imperial policy luxury	16	Metric
Skiing Equipment	Skiing Equipment touring bump stewardship	18	Scalar

Skiing Vacation Deals	Skiing Vacation Deals earth run pass	23	Association
Freestyle Skiing Competitions	Freestyle Skiing Competitions tricia heritage levels	36	Metric
Skiing Accident Prevention	Skiing Accident Prevention organisations turn islands	27	Scalar
Olympic Skiers Profiles	Olympic Skiers Profiles zao run handful	17	Association
Skiing Gear Maintenance	Skiing Gear Maintenance seeing local run	16	Metric
Skiing Levels Progression	Skiing Levels Progression menu run local	28	Scalar
Skiing Events Near Me	Skiing Events Near Me loveyouradk local run	14	Association
Water Skiing Techniques	Water Skiing Techniques menu carp chris	34	Metric
Skiing Cost Factors	Skiing Cost Factors pass produce profit	25	Scalar
Freestyle Skiing Jumps	Freestyle Skiing Jumps run earth stable	12	Association
Skiing Events Today	Skiing Events Today directions local turn	29	Metric
Skiing Cost Estimate	Skiing Cost Estimate run earth policy	11	Association
Freestyle Skiing Tricks	Freestyle Skiing Tricks	31	Metric
Skiing Legends Legacy	Skiing Legends Legacy seeing heritage pioneer	22	Scalar
Skiing Gear Recommendations	Skiing Gear Recommendations compartment like	37	Association

	jokes		
Skiing Levels Assessment	Skiing Levels Assessment run jerry heritage	15	Metric
Water Skiing Locations	Water Skiing Locations fortified chapman au	21	Association
Skiing Cost Calculator	Skiing Cost Calculator you're keep app	18	Metric
Freestyle Skiing Styles	Freestyle Skiing Styles keep thread epic	9	Scalar
Skiing Gear Upgrades	Skiing Gear Upgrades america new prints	33	Association
Skiing Levels Evaluation	Skiing Levels Evaluation Evaluation zermatt procedures program	21	Metric
Skiing Events This Weekend	Skiing Events This Weekend march fall prints	8	Association
Freestyle Skiing Techniques	Freestyle Skiing Techniques	13	Association
Skiing Achievements	Skiing Achievements	17	Metric
Skiing Gear Comparison	Skiing Gear Comparison	27	Scalar
Skiing Levels Progress	Skiing Levels Progress	26	Association
Skiing Events Nearby	Skiing Events Nearby	14	Metric
Water Skiing Basics Guide	Water Skiing Basics Guide	7	Scalar
Skiing Cost Analysis	Skiing Cost Analysis	19	Association
Freestyle Skiing Competitors	Freestyle Skiing Competitors	24	Metric

Observation:

- 1.In most of the queries, Association method seems to give better and more relevant results than scalar and metric clustering.
2. Scalar and metric clustering based query expansion is taking a long time to give the results.

Only few of the local document ids, local stem set and local vocabulary set are provided for each query

Query 1: skiing gear

Local document ids:

b6953acee57e1479bfcc87103dba438
8c7f362735399f5f5e57ed97cf9ed94a
032359deb3293270633f4ede65a65952
67eb00cf5feca9e91e40b6dab2be1f14
8ec257cb0d155724458c59fa38e53ed3
850c0cd1118eb78e16860607466cf4fd
701bae5427754f82f881b5e7f4617274
0f01b02e2aa3ed9c351433a9a24aa53e
c408b6c5b403141a87689feb54c3539e
9c16461a9988d1a8a2118b9522c3e5b9
cab25aa70dde248e9b6589ab6a498212
9f6bb84c769d2ae8ce168a45b9a866c1
df162c4c22c404c90cdf4ddc81b6a6ee
dcfecdc5629798f599fe8079b205c719
81de2f3776c6238527366c44b9ddd51f
eaf970f1c23845c59d07280cd07ef248
fb39d011e21f1278ef377a365c47578a
64cf8b7deac94c31ec06fabdaf551d37
fc3ecf9b8a38638a98978803beb55258
6bc320a35044c47967513647dbde91df
6b64de68adeeb8ac6262699fda7c62a1
524f154895bd4adac911d6d3df77720a

Local stem set:

'hinterthierse' ['hinterthiersee'], 'thierse' ['thiersee'], 'skisport' ['skisport'], 'hartl' ['hartl'],
'rie' ['ried'], 'skiduthyrn' ['skiduthyrning'], 'hyra' ['hyra'], 'skidor' ['skidor'], 'och' ['och'],
'tyrol' ['tyrol'], 'se' ['se'], 'wilder' ['wilder'], 'kaiser' ['kaiser'], 'webbkameror'
['webbkameror'], 'direktrapport' ['direktrapporter'], 'boend' ['boende'], 'planera'
['planera'], 'din' ['din'], 'resa' ['resa'], 'skidkurs' ['skidkurser'], 'vuxna' ['vuxna'], 'barn'

['barn'], 'kontrollera' ['kontrollera'], 'priser' ['priser'], 'med' ['med'], 'billigar' ['billigare'], 'erbjudand' ['erbjudande'], 'fuch' ['fuchs'], 'brixen' ['brixen'], 'thale' ['thale'], 'liftweg' ['liftweg'], 'mobil' ['mobiler'], 'mobile'], 'michel' ['michel'], 'kirchberg' ['kirchberg'], 'deliveri' ['delivery'], 'materi' ['material'], 'accommode' ['accommodation'], 'rainer' ['rainer'], 'scheffau' ['scheffau'], 'hervi' ['hervis'], 'johann' ['johann'], 'rudi' ['rudi'], 'gerri' ['gerry'], 'aschauer' ['aschauer'], 'hochzillert' ['hochzillertal']
 'hollau' ['hollaus'], 'stumm' ['stumm'], 'bike' ['bike'], 'vital' ['vital'], 'kaltenbach' ['kaltenbach'], 'schmiedau' ['schmiedau'], 'stock' ['stock'], 'ww' ['ww'], 'walter' ['walter'], 'maurach' ['maurach'], 'achense' ['achensee'], 'nordic' ['nordic'], 'rofan' ['rofan'], 'achenseestr' ['achenseestr'], 'snow' ['snow']

Local vocabulary set: Displaying few vocabulary

'Photography', 'personalise', 'accompanying', 'material', 'outerwear', 'contain', 'delighted', 'edit', 'advertise', 'form', 'ambassadors', 'probe', 'claude', 'slope', 'es', 'envy', 'support', 'tweet', 'include', 'sites', 'eyeballs', 'history', 'waterville', 'alto', 'ask', 'especially', 'struggle', 'sportsways', 'river', 'take', 'born', 'pro', 'completely', 'live', 'skimming', 'women', 'placid', 'cat', 'preferences', 'okay', 'wordpress', 'runtothefinish', 'incredibly', 'oct', 'hut', 'stuck', 'means', 'leases', 'orbital', 'slang', 'reports', 'impact', 'online', 'speedy', 'display', 'decide', 'verbal', 'machine', 'podcasts', 'value', 'visit',

'ballclapper', 'priceless', 'goggles', 'search', 'encouraging', 'deal', 'renter', 'patrick', 'service', 'blister', 'instead', 'markets', 'black', 'essentials', 'mogul', 'stickers', 'telluride', 'refine', 'youd', 'know', 'developed', 'charged', 'pigon_racer', 'open', 'sitemap', 'offers', 'right', 'update', 'sure', 'edited', 'people', 'self', 'company', 'fresh', 'change', 'painful', 'accounting', 'ahead', 'channel', 'wasnt', 'reply', 'shops', 'crested', 'towards', 'poster', 'theyd', 'tips', 'independent', 'prix', 'pay', 'phone', 'act', 'knows', 'audience', 'year', 'avoid', 'cost', 'talking', 'merino', 'kid', 'language', 'commission', 'fashion', 'barnes', 'town', 'javascript', 'burton', 'mike'

Local stem set:

'zobraz' ['zobraz'], 'easter' ['easter'], 'holiday' ['holiday'], 'holidays'], 'guildfordskislopecouk' ['guildfordskislopecouk'], 'skip' ['skip'], 'content' ['content'], 'slope' ['slopes'], 'slope'], 'donut' ['donutting'], 'parti' ['parties'], 'info' ['info'], 'find' ['find'], 'due' ['due'], 'staff' ['staff'], 'close' ['closed'], 'close'], 'monday' ['monday'], 'april' ['april'], 'th' ['th'], 'open' ['opening'], 'open'], 'tuesday' ['tuesday'], 'current' ['currently'], 'current'], 'remaind' ['remainder'], 'period' ['period'], 'ndrd' ['ndrd'], 'fulli' ['fully'], 'apart' ['apartments'], 'apart'], 'public' ['public'], 'session' ['sessions'], 'daytim' ['daytime'], 'even' ['even'], 'evenings', 'evening'], 'morn' ['morning'], 'pm' ['pm'], 'bore' ['boring'], 'stuff' ['stuff'], 'must' ['must'], 'advanc' ['advance'], 'enquir' ['enquire'], 'pleas' ['please'], 'enquiri' ['enquiry'], 'form' ['form'], 'karin' ['karin'], 'bradburi' ['bradbury'], 'share' ['share'], 'choos' ['choose'], 'platform' ['platform'], 'facebook' ['facebook'], 'x' ['x'], 'reddit' ['reddit'], 'linkedin' ['linkedin'], 'tumblr' ['tumblr'], 'pinterest' ['pinterest'], 'vk' ['vk'], 'email' ['email'], 'copyright' ['copyright'], 'christ' ['christians'], 'colleg' ['college'], 'club' ['club'], 'sitemap' ['sitemap'], 'legal' ['legal'], 'page' ['pages'], 'page'], 'load' ['load'], 'link'

['link'], 'go' ['go', 'going'], 'octob' ['october'], 'half' ['half'], 'look' ['looking'], 'someth' ['something'], 'bring' ['bring'], 'privat' ['private'], 'week' ['week'], 'upward' ['upwards'], 'practic' ['practice'], 'instructor' ['instructors', 'instructor'], 'given' ['given'], 'one' ['one'], 'passport' ['passports', 'passport'], 'give' ['give'], 'call' ['call'], 'minimum' ['minimum'], 'age' ['age'], 'accompani' ['accompanied'], 'per' ['per'], 'person' ['person'], 'everi' ['every'], 'time' ['times', 'time']

Query 2: skiing resort

Local Document Set:

81de2f3776c6238527366c44b9ddd51f
eaf970f1c23845c59d07280cd07ef248
fb39d011e21f1278ef377a365c47578a
64cf8b7deac94c31ec06fabdaf551d37
fc3ecf9b8a38638a98978803beb55258
6bc320a35044c47967513647dbde91df
6b64de68adeeb8ac6262699fda7c62a1
524f154895bd4adac911d6d3df77720a
60f4df1a17154e65dd028570f02617a0
2be3849678b7c891641004fcf84ed051
ebf323c219ea9edb62951b0d217ed507
9fde3613023f0ec5019827226d6c78cf
4d04ae853da46b1f0c26d396b0d1eff8
67b6e0d88c453ef0919a4182b2ad5074
648a0c52e4b47c2e4547f7f8dc4ec405
d6174b4c0433ab21e1bf48263a3b18ae
db6953acee57e1479bfcc87103dba438
8c7f362735399f5f5e57ed97cf9ed94a
032359deb3293270633f4ede65a65952
67eb00cf5fec9e91e40b6dab2be1f14
8ec257cb0d155724458c59fa38e53ed3
850c0cd1118eb78e16860607466cf4fd
701bae5427754f82f881b5e7f4617274
0f01b02e2aa3ed9c351433a9a24aa53e
c408b6c5b403141a87689feb54c3539e
9c16461a9988d1a8a2118b9522c3e5b9
cab25aa70dde248e9b6589ab6a498212
9f6bb84c769d2ae8ce168a45b9a866c1

Local vocabulary set: Displaying few vocabulary

'mobile', 'page', 'era', 'collecting', 'ballet', 'rebell', 'becoming', 'lightweight', 'posted', 'promotion', 'oldmanski', 'drop', 'skip', 'team', 'wherewhat', 'país', 'resort', 'passion', 'muscle', 'joy', 'wild', 'majics', 'innovative', 'freakschoolers', 'email', 'mesa', 'plastic', 'guests', 'nah', 'literally', 'guides', 'help', 'shot', 'spect', 'bumps', 'edges', 'anyone', 'magnetise', 'narrow', 'joal', 'inspiredmedia', 'rainbow', 'saw', 'government', 'guest', 'holiday', 'waterproof', 'enough', 'notorious', 'age', 'sometimes', 'knowledge', 'resident', 'skissnowboards', 'attend', 'sun', 'heuga', 'chart', 'wanted', 'commissions', 'dsa_departmentwinter', 'message', 'safety', 'grits', 'news', 'journal', 'routes', 'helitrap', 'atomic', 'say', 'discounts', 'climbing', 'insider', 'dash', 'muffman', 'iceland', 'jan', 'nonetheless', 'already', 'foremost', 'recent', 'side', 'purchased', 'enable', 'pond', 'norway', 'charger', 'lift', 'save', 'tent', 'chance', 'weekend', 'seems', 'nervous', 'great', 'cookies', 'heres', 'adrenaline', 'greater', 'someone', 'boot', 'theres', 'thought', 'covet', 'bookstores', 'beanie', 'skitalk', 'plake', 'rentals', 'starglass', 'shutterstock', 'uk', 'bigpurpleskisuit', 'versus', 'rider', 'category', 'founder', 'therefore', 'exchanges', 'coloradodogfart', 'ease', 'memory', 'using', 'archive', 'containing', 'chunderface', 'httpswwwpugskicomforumsheritage', 'cooper', 'hassle', 'jacket', 'ive', 'dances', 'czechoslovakia', 'original', 'feet', 'bikes', 'posters', 'cdubdashrdr', 'furniture', 'norsquip', 'ultimate', 'a_platic_bag', 'hiking', 'web', 'karen_skier', 'independent', 'prix', 'pay', 'phone', 'act', 'knows', 'audience', 'year', 'avoid', 'cost', 'talking', 'merino'

Local stem set:

'slot' ['slots'], 'Also' ['also'], 'Offic' ['office'], 'hour' ['hours'], 'friday' ['friday'], 'want' ['want'], 'obvious' ['obviously'], 'anytim' ['anytime'], 'like' ['like'], 'swisspassescom' ['swisspassescoms', 'swisspassescom'], 'swiss' ['swiss'], 'rentalhir' ['rentalhire'], 'resort' ['resorts', 'resort']

Query 3:ski board**Local Document set:**

81de2f3776c6238527366c44b9ddd51f
eaf970f1c23845c59d07280cd07ef248
fb39d011e21f1278ef377a365c47578a
64cf8b7deac94c31ec06fabdaf551d37
fc3ecf9b8a38638a98978803beb55258
6bc320a35044c47967513647dbde91df
6b64de68adeeb8ac6262699fda7c62a1
524f154895bd4adac911d6d3df77720a
60f4df1a17154e65dd028570f02617a0

2be3849678b7c891641004fcf84ed051
 ebf323c219ea9edb62951b0d217ed507
 9fde3613023f0ec5019827226d6c78cf
 4d04ae853da46b1f0c26d396b0d1eff8
 67b6e0d88c453ef0919a4182b2ad5074
 648a0c52e4b47c2e4547f7f8dc4ec405
 d6174b4c0433ab21e1bf48263a3b18ae
 db6953acee57e1479bfcc87103dba438
 8c7f362735399f5f5e57ed97cf9ed94a
 032359deb3293270633f4ede65a65952
 67eb00cf5fec9e91e40b6dab2be1f14
 8ec257cb0d155724458c59fa38e53ed3
 850c0cd1118eb78e16860607466cf4fd
 701bae5427754f82f881b5e7f4617274
 0f01b02e2aa3ed9c351433a9a24aa53e
 c408b6c5b403141a87689feb54c3539e
 9c16461a9988d1a8a2118b9522c3e5b9
 cab25aa70dde248e9b6589ab6a498212
 9f6bb84c769d2ae8ce168a45b9a866c1

Local vocabulary set: Displaying few vocabulary

'permit', 'lyžeřské', 'nabídku', 'checked', 'granettes', 'liftweg', 'family', 'pm', 'planuj',
 'kelchsau', 'se', 'really', 'shop', 'latest', 'neben', 'juwel', 'join', 'opening',
 'taniacourchevel', 'da', 'och', 'dospěl', 'mostly', 'weekdays', 'mit', 'unnecessary',
 'excursions', 'skischulen', 'whilst', 'děti', 'certain', 'planner', 'magazin', 'überblick', 'dorf',
 'skis', 'skipasú', 'money', 'til', 'sk', 'story', 'taken', 'count', 'pass', 'working', 'driving',
 'week', 'urban', 'österrike', 'parties', 'remainder', 'tyrolsko', 'sturtevants', 'view', 'load',
 'newslettersblogs', 'lessons', 'gallery', 'number', 'poll', 'sport', 'skileje', 'getting',
 'anytime', 'medical', 'thiersee', 'ziller', 'sv', 'time', 'swiss', 'last', 'kinder', 'flight',
 'combine', 'valley', 'viaje', 'financially', 'hours', 'karnetów', 'og', 'march', 'ie', 'aid',
 'guildford', 'oferta', 'tras', 'førstehåndsrapporter', 'packages', 'slots', 'karin', 'menu',
 'priser', 'hervis', 'snörapport', 'things', 'connys', 'søk', 'fully', 'month', 'order', 'der', 'fly',
 'rights', 'markbachjochbahn', 'wont', 'fr', 'apartments', 'taniej', 'events', 'forest', 'vær',
 'like', 'details', 'overnatting', 'closed', 'tirol', 'term', 'locaties', 'vital', 'early', 'das', 'book',
 'château', 'whats', 'bundesstrasse', 'hotel', 'nearby', 'adults', 'well', 'cz', 'schmiedau',
 'seasonaires', 'maps', 'informes', 'soleil', 'exposure', 'season', 'hollaus', 'berglift', 'av',
 'contact', 'half', 'schneebericht', 'otp', 'possibility', 'coach', 'thyon', 'raporty', 'powder',
 'sprawdz', 'schneeberichte', 'due', 'alpinschiverleih', 'need', 'though', 'pertisau', 'roku',
 'plan', 'trati', 'kleinriedstra e', 'every', 'pricing', 'mosses', 'main', 'allow',
 'wypo yczalnie', 'eases', 'go', 'precios', 'zpr avy', 'leave'

Local stem set:

'report' ['reports', 'report'], 'plan' ['plan', 'planning'], 'relat' ['related'], 'nearbi' ['nearby'], 'weather' ['weather'], 'first' ['first'], 'hand' ['hand'], 'overview' ['overview'], 'review' ['reviews'], 'pist' ['piste'], 'map' ['map', 'maps'], 'planner' ['planner'], 'lift' ['lifts', 'lift'], 'ticket' ['tickets'], 'children' ['children'], 'year' ['year', 'years'], 'search' ['search'], 'check' ['check', 'checked'], 'price' ['prices', 'pricing'], 'avail' ['availability', 'available'], 'view' ['view'], 'deal' ['deal'], 'lesson' ['lessons'], 'skischulen' ['skischulen'], 'materialzustellung' ['materialzustellung'], 'waar' ['waar'], 'kun' ['kun'], 'je' ['je'], 'skimateria' ['skimateriaal'], 'huren' ['huren'], 'steinberg' ['steinberg'], 'een' ['een'], 'overzicht' ['overzicht'], 'van' ['van'], 'skiverhuur' ['skiverhuur'], 'locati' ['locaties'], 'sneeuwrapport' ['sneeuwrapport'], 'gelinkt' ['gelinkt'], 'aan' ['aan'], 'achenseetirol' ['achenseetirol'], 'oostenrijk' ['oostenrijk'], 'dichtbij' ['dichtbij'], 'christlum' ['christlum'], 'pertisau' ['pertisau'], 'weer' ['weer'], 'gebruikersrapport' ['gebruikersrapport'], 'verblijf' ['verblijf'], 'jouw' ['jouw'], 'accommodato' ['accommodatoe'], 'volwassen' ['volwassen'], 'kinderen' ['kinderen'], 'onder' ['onder'], 'jaar' ['jaar'], 'zoek' ['zoek'], 'control' ['controleer'], 'prijsen' ['prijsen'], 'en' ['en'], 'beschikbaarheid' ['beschikbaarheid'], 'bij' ['bij'], 'onz' ['onze'], 'goedkop' ['goedkoper'], 'bekijk' ['bekijk'], 'aanbiedingen' ['aanbiedingen'], 'skiutlei' ['skiutleie'], 'leie' ['leie'], 'av' ['av'], 'langrenn' ['langrenn'], 'og' ['og'], 'skiinfono' ['skiinfono'], 'tilknyttet' ['tilknyttet'], 'webkamera' ['webkamera'], 'overnat' ['overnatting'], 'planlegg' ['planlegg'], 'tur' ['tur'], 'innkvart' ['innkvartering'], 'voksn' ['voksne'], 'sjekk' ['sjekk'], 'billiger' ['billigere'], 'tilbud' ['tilbud'], 'skitrim' ['skitimer'], 'nart' ['nart'], 'pogoda' ['pogoda'], 'planuj' ['planuj'], 'wyjazd' ['wyjazd'], 'w' ['w'], 'kameri' ['kamery'], 'raporti' ['raporty'], 'ze' ['ze'], 'stoku' ['stoku'], 'opini' ['opinie'], 'mapa' ['mapa'], 'tra' ['tras'], 'opadi' ['opady'], 'noclegi' ['noclegi'], 'ceni' ['ceny'], 'zoplanuj' ['zoplanuj'], 'zakwaterowani' ['zakwaterowanie'], 'osobi' ['osoby'], 'dzieci' ['dzieci'], 'roku' ['roku'], 'szukaj' ['szukaj'], 'u' ['u'], 'naszego' ['naszego'], 'partnera' ['partnera'], 'taniej' ['taniej'], 'zobacz' ['zobacz'], 'narciarski' ['narciarskie'], 'dostawa' ['dostawa'], 'miejscy' ['miejscy'], 'zakwaterowania' ['zakwaterowania'], 'alquil' ['alquiler'], 'equipo' ['equipos'], 'niev' ['nieve'], 'relacionada' ['relacionadas'], 'cercano' ['cercano'], 'tiempo' ['tiempo'], 'inform' ['information'], 'informes'], 'alojamiento' ['alojamiento'], 'planifica' ['planifica'], 'tu' ['tu'], 'viaj' ['viaje'], 'adulto' ['adultos'], 'menor' ['menores'], 'buscar' ['buscar'], 'consult' ['consulte'], 'precio' ['precios'], 'disponibilidad' ['disponibilidad'], 'con' ['con']

Collaboration with the UI:

1. The query is received from the UI through an API call.
2. We receive the query and the query expansion method chosen by the user as the parameters to the query expansion part of the code. Then the documents(Solr results) are retrieved from the Solr collection. These documents are considered as local documents, and local vocabulary sets and local stem sets are created.

3. Based on the query expansion method(Association, Scalar or Metric) specified by the user an expanded query is formed and this query is fed to the relevance models
4. The relevance model gets the result for the expanded query from the Solr indexed collection which is then passed to the UI to be visible to the user.

7. Discussion

Upon completing our project, we've acquired a comprehensive understanding of the multifaceted elements necessary for a functional search engine. Our efforts resulted in a basic search tool capable of providing somewhat relevant results by applying the techniques taught in class. This endeavor not only provided us with practical insights into search engine development but also underscored the vastness of the field of Information Retrieval. A comparison of our outcomes with the prowess of major search engines like Google and Bing revealed the considerable gap between our efforts and the complexities involved in maintaining such widely-used platforms. It's evident that our project only scratched the surface of what's needed for a successful and current search engine.

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 Religions. The five of us split 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 Skiing. 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.