# Project 2

# Web Search Engine – Travelpedia

## COEN 272 – Web Search and Information Retrieval
Department of Computer Science
Santa Clara University

| Kavya Gautham | Madhumitha Mani | Nisha Narayanaswamy | Priyanka Rao |
|---|---|---|---|
| W1166971 | W1157297 | W1153362 | W1166161 |
| kgautam@scu.edu | mmani@scu.edu | nnarayanaswamy@scu.edu | prao@scu.edu |

### Abstract
The project aims at developing a robust web search engine for the travel domain. The targeted user base is anyone who is interested in finding vacation destinations within California. The search engine is built on ElasticSearch leveraging its various features. The goal of this project is to provide relevant search results to the users through various techniques like query boosting, query expansion, stemming, stopping, spell check, query parsing, handling geo-distance queries and automatic facet selection, which are in addition to the default features provided by Elasticsearch. The success of the web search engines: traditional and faceted, is determined by the effectiveness, efficiency and user satisfaction derived from the user evaluation studies.

## 1 Search Engine Overview

### 1.1 Introduction

The project is divided into five parts: (i) Crawl the travel related websites in California (ii) Build a traditional search engine using Elastic Search (iii) A non-traditional interaction feature using faceted search interface (iv) Logging component for user queries to log queries, clicks, session duration (v) Evaluation of the basic and advanced versions of the search engine through a user study. The architecture and the various features of the web search engine (Travelpedia) are detailed in the following sections.

### 1.2 High-level Architecture

The high-level architecture of the search is shown in the Figure 1.2. There are four stages in the development of Travelpedia – travel domain search engine (i) Crawling – getting the webpages related to the travel domain (ii) classifying/automatic labeling of the documents (iii) Indexing – involves creating mappers and analyzers for the documents (iv) Retrieval – process of retrieving documents from the search engine
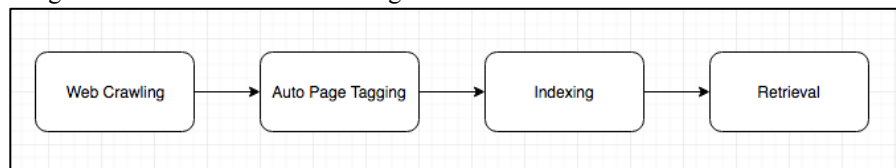


Figure 1.2 Search engine high-level architecture

## 2 Crawling and Document Indexing

### 2.1 Automatic document labeling
The automatic document labeling involves the following steps:
- Web pages from specific popular cities within California have been chosen and shortlisted for crawling. For example, the travel giants like travelpedia, visitcalifornia, discovercalifornia and many more specific sites for destinations throughout California were considered as the main seeds.

- A set of 400 specific sites covering all the facet categories were first chosen, crawled and manually labeled. Next, a set of 1300 more travel related pages were crawled and extracted for Auto Labeling process. The Auto Labeling process is described as follows:
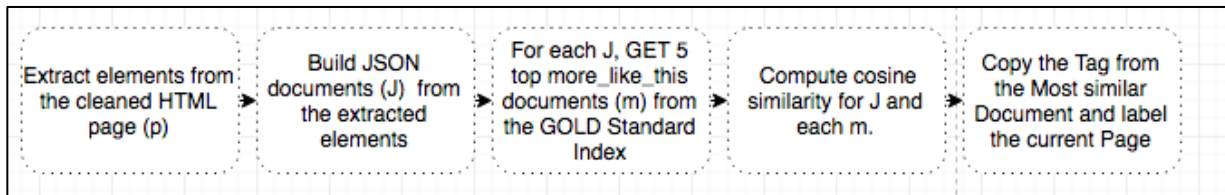


Figure 2.1 Automatic document labeling

## 2.2    Document indexing

As part of pre-processing we do the following: when each HTML file is read, unwanted spaces and empty lines are removed. Then, the file is parsed through JSoup, to create the DOM structure. Once the DOM structure is created, script, comment, header, br, style, meta, CDATA tags are removed using JSoup, with the assumption that these tags will not provide any relevant information.  Figure 2.2 shows the various steps involved before a document gets indexed into ElasticSearch.
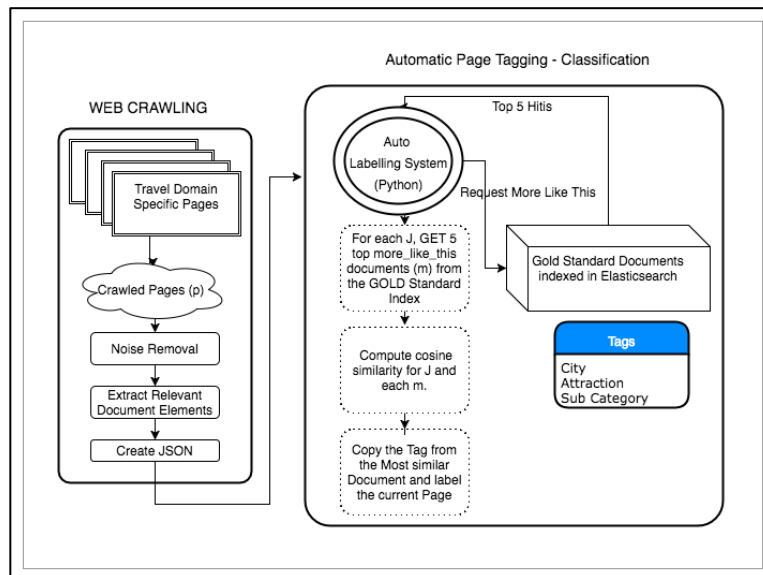


Figure 2.2 Document indexing

The technique is based on the assumption all the viewable parts in a webpage are contained in the <body> tag. All the new metrics calculated were added as attributes to the nodes hence, keeping track of the metrics and making decisions per-node level was easy.

Before the documents are indexed, custom analyzers and mapping were created.  Custom analyzers were created for stemming (Porter stemmer), stopping (English stopwords removed), lowercase, trimming (to remove extra spaces). The mapping uses the created custom analyzers for the various fields accordingly.

## 3    Document retrieval architecture

The search engine Travelpedia is built using Play framework. Play framework makes it easy to build web-applications with Java provides a lightweight, stateless, web-friendly architecture.  We chose play framework as it provides a holistic platform for integrating the different components of the search engine like UI, logic processing component, interface to Elasticsearch. The following are the major components of the document retrieval model using Play framework:

## 3.1    View

The view is responsible for generating the HTML pages for user interaction with the search engine. There are two types of html pages generated: (i) Traditional Search Engine UI (ii) Faceted Search Engine UI. Both views provide user interface for login/logout component, executing search queries, viewing the SERP and pagination. In addition to these interfaces, faceted search engine UI has an additional component that enables users to apply filters for queries. Any action on the view elements triggers the execution of corresponding action in the controller.
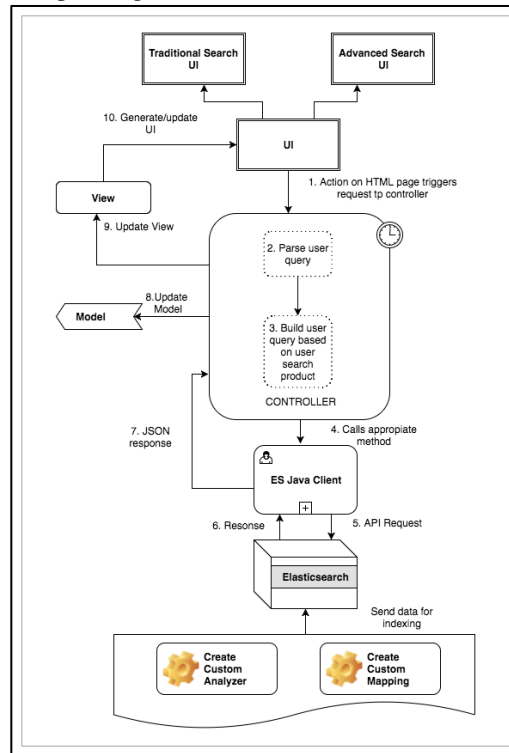


Figure 3.1 Document retrieval architecture

## 3.2    Controller

The controller is the crucial component of the document retrieval architecture as it wires together different components of the search engine. The controller handles all requests from the view page elements and populates the Model components necessary to perform search. Once all the relevant information is gathered, it invokes the corresponding service in the Java Client. Controller also waits to receive the response from the Java client which is in a JSON format, parse the JSON object to retrieve relevant results and populate the Model components that are necessary to render the SERP back to the user.

## 3.3    Java Client – Java interface to Elasticsearch

Elasticsearch is programming language independent and the APIs required for indexing, searching etc. could be accessed using HTTP and JSON in any language. Since the Travelpedia development team has strong Java expertise, it was a natural choice to use ElasticSearch Java client. The java client is created as a Transport client to perform API calls to ElasticSearch. It receives requests from the controller and converts it into JSON queries that are executed at the Elasticsearch. The response from the ElasticSearch is also in the JSON format, which is sent the controller that initiated this request.

## 3.4    Model

The model maintains the data needed to be passed across the requests and responses to and fro from UI and ElasticSearch. For example: the data retrieved from the search engine is stored as "SearchResult" objects, which contains a list of "Place" model objects. The Place model object is used to store search results returned from Elasticsearch like document title, liveURL, snippet and city information. Once the Controller populates the model objects, it also updates the View with the relevant information.

# 4        Versions of Search Engine

As mentioned earlier, the project demanded the development of two versions of search engines: the traditional and the advanced version. The underlying indexing and document retrieval architecture remains the same for both the search engines with just one striking difference in the advanced version, which is dealt in the following section.

## 4.1        Traditional Search Engine

The traditional search engine provides users with a simple interface to query documents related to travel destinations. The elements of the traditional SERP are **Input Box** - place to enter query, **Sign in button** - to login to the system, needed for user log, **Query correction** - replaces spelling error and provide suggestion, **Result count** - number of documents found for a particular query, **Title** - Document name, **Link** - url of the document, **Snippet** - small paragraph describing the document with query
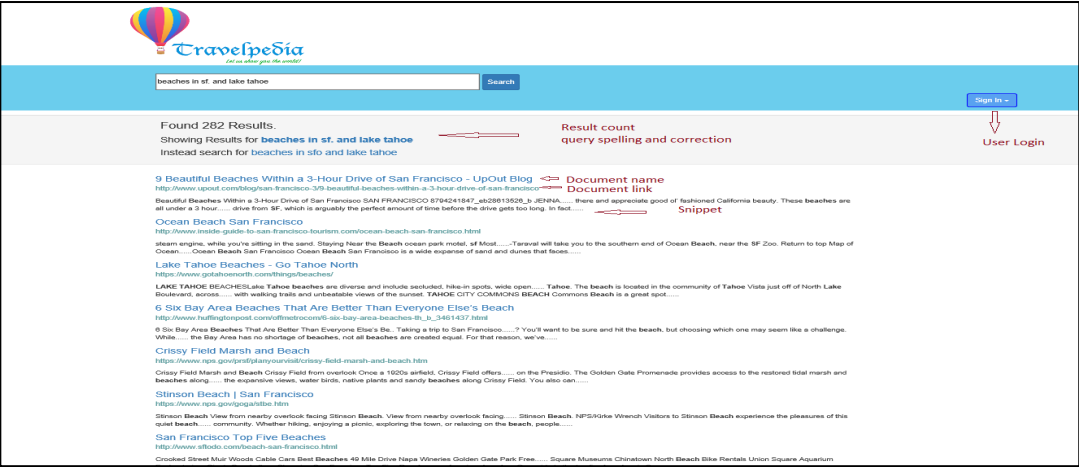


Figure 4.1 Traditional SERP

word highlighted, **Pagination** - Result is divided into pages and each page displays 11 results. User can click next page number to view next page.

## 4.2        Faceted Search Engine

For the advanced version of the search engine, we decide to create facets for the following reasons: (i) there is an inherent hierarchy for the destinations like Parks which can be further classified as national/state park similarly Museums classified as Art, Science and History (ii) allows users to navigate easily according to several different topics than in a single, pre-determined order (iii) help users to learn new subcategories of the search. Figure 4.2 shows all the facets listed in our search engine.
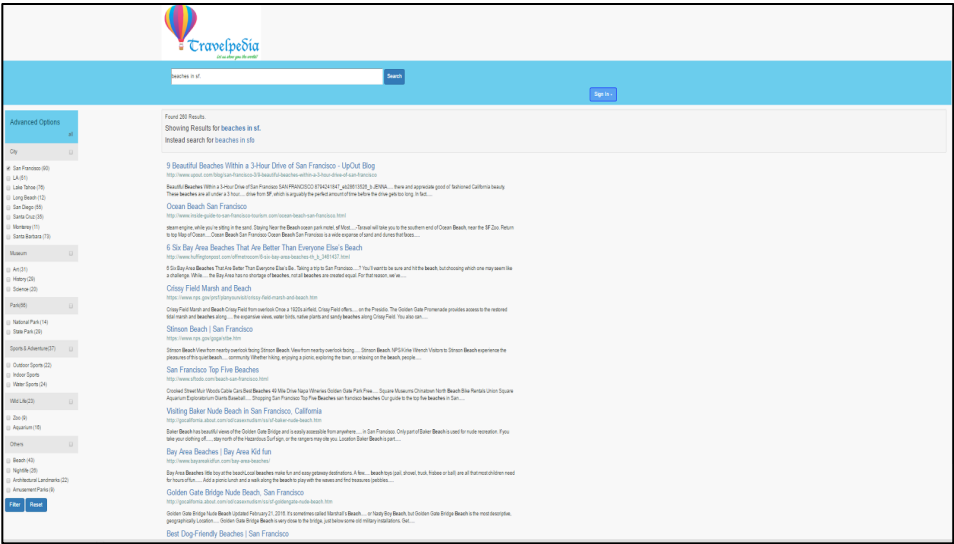


Figure 4.2 Faceted SERP

User can select one facet at a time or can select multiple facets using check-all checkbox, which selects group of facets. User can view the results of faceted query after clicking filter button. User can use reset button to unselect the facets selected in earlier query. There is no limit on facet selection. The results of the faceted query are displayed on right side of the same page. Each time a query is entered, the document counts of all facets are retrieved and displayed in facet drop down box.

## 5. General search engine features

There are a number of features that are common to both search engines as mentioned in the following section. The features were kept common as a design choice because we did not want to introduce bias towards any system.

### 5.1 Query-time Boosting

Boosting is an important tool used to tune relevance. Elastic search offers two types of boosting: (i) Index boosting (ii) Query-time boosting. In our case, index boosting was not necessary, as we had created only one index for all Travelpedia related documents. Figure 5.1 shows the snapshot of the query with boosting. Practically, there is no simple formula for specifying the "perfect" boost value for a query clause. It is a try-it-and-see task and finally depends on the design choice of the developers.

```
"query": {                                    ▶ 🔧
    "bool": {
        "should": [
            {
                "match": {
                    "city": {
                        "query": "things to do in san francisco",
                        "analyzer": "query_analyzer",
                        "boost": 50
                    }
                }
            },
            {
                "match": {
                    "attraction": {
                        "query": "things to do in san francisco",
                        "boost": 20
                    }
                }
            },
            {
                "match": {
                    "title": {
                        "query": "things to do in san francisco",
                        "boost": 10
                    }
                }
            },
            {
                "match": {
                    "content": {
                        "query": "things to do in san francisco",
                        "boost": 0.5
                    }
                }
            },
```

Figure 5.1 Query with boosting

ElasticSearch inherently supports "natural" boosting based on the "field-length norm". In Travelpedia, there was clearly a difference in the importance over various fields. For example: if the query terms matched against a city, it was boosted by a factor of 50, attraction and sub-category by a factor of 20 followed by title boosted by a factor of 10, keywords by 5 and finally the content by 0.5. This combination of boost values helped us achieve an acceptable level of document relevance across different query varieties.

### 5.2 Query Expansion

In most documents, the same concept may be referred using different words. In Travelpedia, for example: A user might search for "beaches in LA" and expect to find documents that contain "los angeles, la, losangeles, L.A." This issue is known as "synonymy" and has an impact on the overall recall of the information retrieval system. We overcome this problem only at the "city" level by introducing "synonyms" for the various cities we have at the time of indexing.

Figure 5.2 Query Analyzer with synonym for cities

As shown in Figure 4.2, query expansion is achieved using the concepts of "synonyms" in ElasticSearch. This helps expand the search query containing a city, into a number of queries with city variations based on the synonyms defined at the time of document indexing. This query expansion logic happens under the hood transparent to the user.

## 5.3    Snippet Generation and keyword highlight

Snippets are an important element of the SERP because it helps user get a sneak-peek of the document contents and help determine document relevance. Even more sophisticated snippets show "highlights" of text from the search result so a user can easily understand why the document matched the query.
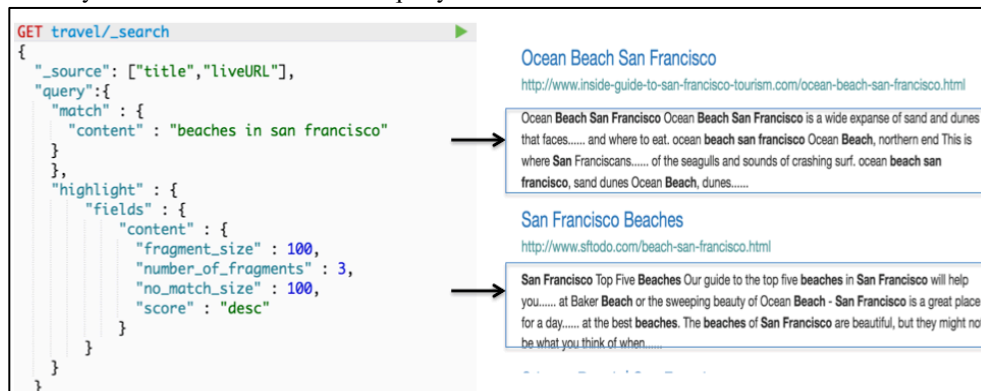


Figure 5.3 Highlight query with generated snippets

Travelpedia achieves snippet generation with keyword highlighted using the "highlight" parameter in the search query. The highlight query code fragment and an example snippet are shown in Figure 5.3. The snippets are generated as fragments with a maximum of 3 fragments per snippet. Matching query keywords against the content field generates the snippet. Elasticsearch internally scores the fragments, and we use the top 3 fragments for display and to highlight the corresponding query keyword. If the keyword is not present in the content, first 100 words in the content will be printed.

## 5.4    Query Parsing to support geo-distance queries

When a user submits the search query, Travelpedia parses the query before any results are fetched from Elasticsearch. The aim of the query parser is to identity if the user has specified  "distance-based" queries. For example: if the user is looking for state parks within 200 miles, the query parser identifies the distance as "200 miles" and narrows the search to places that are within the specified distance. The query and the corresponding results are shown in Figure 5.4

```
GET /travel/_search                              ▶  ⚙
{
    "query": {
        "filtered": {
            "filter": {
                "geo_distance": {
                    "distance": "200 miles",
                    "coordinates": {
                        "lat":  40.715,
                        "lon": -73.988
                    }
                }
            }
        }
    }
}
```

Figure 5.4 Geo-distance query

If the query recognizes it as a distance-based query, the search is restricted to all the places that are within the specified distance. The user's IP address is captured and converted into geo-co-ordinates (latitude and longitude). The geo-coordinates are then added to the ElasticSearch query to get the appropriate results.

## 5.5     Spell check option

Similar to Google spell correction, Travelpedia also check for any spell-check errors in the user's search query. This feature leverages the "term" suggester API, which suggests terms based on edit distance.

The suggestions are given for those terms that are missing from the "content" field. For the purpose of spell check, we create another version of the content field where the data is not analyzed/stemmed/stopped. Travelpedia considers the top 1 correction suggested by Elasticsearch. One constraint for spell check is that the input data to term suggester should be a minimum of 3 characters in length. Figure 5.5 shows the example of the term suggester in action.



Figure 5.5 Spell check example

## 5.6     Pre-selection of city facet

In the faceted search engine, if the user types cityname (which happens to one of the facets) in the search query, Travelpedia identifies it and pre-selects the city facet for the user. This is achieved by parsing the query, looking for possible city candidates that might be mentioned in the query. If the query parser returns any result, then we pre-select all the cities that the user has mentioned in the query.
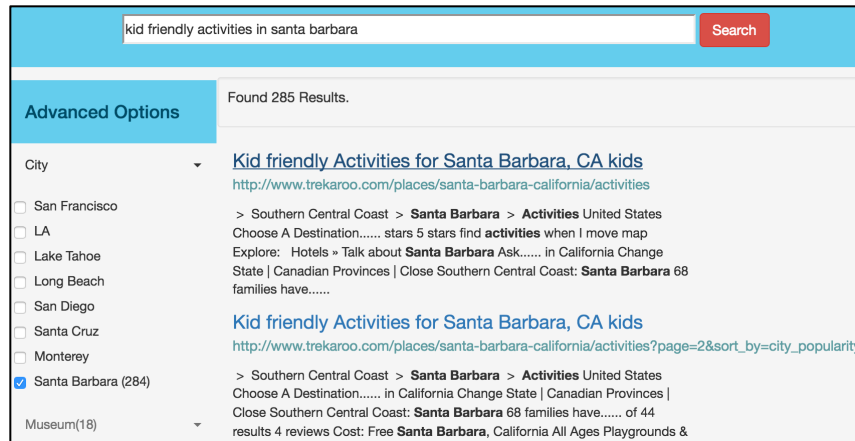
Figure 5.6 Pre-selection of city facet

# 6    User Logging

For the purpose of search engine evaluation, we created logs of all the user activities across the search engines. We have created 5 user credentials for both traditional and faceted system. We are maintaining 5 user log files one for each user.  Each user interaction is logged into his respective file.

## 6.1    Need for user logging

User logging helps us to understand the system from user perspective.  For example: session time variable gives us an idea of time spent by each user on our system.  The query list and pages viewed gives us an idea of which documents a user found more relevant for particular query and how many pages he had to click until he found his relevant document. These components can be used in future to improve the search engine. For example if many users viewed a subset of document for a particular query, we can boost such document up in the index so that next time these pages come as top results. User logging also helps in providing suggestion and recommendation to user.

## 6.2    Elements of the logs

Each time a user logs into the system we keep track of following components:

- Session variables (**start session time** - time at which user logged into the system; **end session time** - time at which user logged out)
- Query List - list of queries searched by a user in that session



> Search Engine Type: Faceted SERP
>
> Session start time: 2016/05/31 12:33:37
>
> Queries: [boat lake tahoe, {cities: [] attraction: [Park] subcategory: [National Park, State Park]},
>
> {cities: [] attraction: [Park] subcategory: [National Park, State Park]},
>
> {cities: [] attraction: [Park] subcategory: [National Park, State Park]},
>
> {cities: [San Diego] attraction: [] subcategory: []} ]
>
> Mouse click count: 8
>
> Pages_viewed : [1, 1, 3, 1, 2]
>
> Session end time: 2016/05/31 12:37:32

Figure 6.1 User Logs

- Mouse clicks  - number of links clicked for each user task
- Pages viewed - gives list of pages viewed for a particular query
- Engine Type - Traditional SERP or Faceted SERP
- For Faceted Search Engine: In addition to above variables we also keep track facet query search. Facet query search contains multiple facets selected by user to filter the result

**6.      User Evaluation**
**6.1      Evaluation goal**

A small user study was conducted to compare the traditional and faceted version of the search engine. As part of usability research once both the search engine versions were developed, activities that a typical user would perform to achieve their goal on travel websites were listed. Five target participants between the ages 20-30 we identified. All were intermediate level users, including students and working professionals.

**6.2      User study methodology**

The procedures followed during the user study are as follows:
- Five target participants between the ages 20-30 we identified.
- Five tasks were identified to execute in both search engines.
- Task allocation matrix was created. Half the users were shown the traditional search engine to evaluate first followed by the faceted search engine. The other half was shown the faceted search first.
- A single system with Mac OSX was chosen for the users to evaluate on.
- Moderator provided initial introduction about the two search engine interfaces to users.
- Each user was given one task on each search engine interface. One user was evaluated at a time.
- Users were asked to login to the system. A new session was started before every task.
- Maximum time limit per task set as 10 minutes.
- Screen recording for all user evaluations, task duration was recorded using a digital stopwatch. The user logs were also recorded in each version of the search engine.
- Moderator did not provide any assistance to the users once the task evaluation started.
- Errors if any during evaluation were recorded.

A usability metric reveals some aspect of the user's interaction with the system. The metrics used for user evaluation were three-fold:
- **Efficiency** is the amount of effort required to complete a task. Measured the time taken, number of clicks, number of queries required completing each task. Since for each search engine, different users executed tasks 1 through 5, an average of each of the metrics across the two versions was calculated.
- **Effectiveness** is being able to complete a task. A task is successful if relevant documents were found accordingly to the user query. Three non-binary task success levels were defined and the user was asked to rate each task as "Success", "Partial Success" or "Failure". No assistance was provided to the user during task evaluation.
- **User satisfaction** determines how happy the user is with a given system. The USE questionnaire was provided to the user after a task was completed in each search interface. A five-point Likert scale ranging from "Strongly Disagree" to "Strongly Agree" was used for rating.

**6.3 User study results**

The following five tasks were executed in both search engines in random order:

| | |
|---|---|
| Task 1 | You are visiting Monterey and would like to know the things you can do there. Try to narrow down to at least 3 activities of your choice. |
| Task 2 | You are planning a bachelor/bachelorette party for your best friend in Los Angeles. You would like select a bar, restaurant for dinner. |
| Task 3 | You want to plan a weekend trip to parks, museums and/or beaches within 200 miles. |
| Task 4 | You would like to go whale watching. Pick a location of your choice, the best time and season to go. |
| Task 5 | Plan your day in San Francisco |

**Figure 6.1: Tasks to be executed in traditional and faceted search engines**

|  | TaskSet | TaskSet |
|---|---|---|
| **User 1** | SE1-Task1 | SE2-Task2 |
| **User 2** | SE2-Task3 | SE1-Task2 |
| **User 3** | SE1-Task3 | SE2-Task4 |
| **User 4** | SE2-Task5 | SE1-Task4 |
| **User 5** | SE1-Task5 | SE2-Task1 |

**Figure 6.2: Tasks allocation matrix**

| User # | Time per task (in mins) | | Number of queries per task | | Number of clicks per task | |
|---|---|---|---|---|---|---|
|  | SE1 | SE2 | SE1 | SE2 | SE1 | SE2 |
| User 1 | 3.26 | 3.4 | 1 | 3 | 4 | 4 |
| User 2 | 1.18 | 1.15 | 1 | 2 | 3 | 1 |
| User 3 | 3.33 | 2.26 | 3 | 2 | 4 | 1 |
| User 4 | 2.15 | 3.15 | 2 | 3 | 3 | 4 |
| User 5 | 4.03 | 4.27 | 3 | 3 | 2 | 5 |
| **Average** | **2.79** | **2.846** | **2** | **2.6** | **3.2** | **3** |

**Figure 6.2: Efficiency**

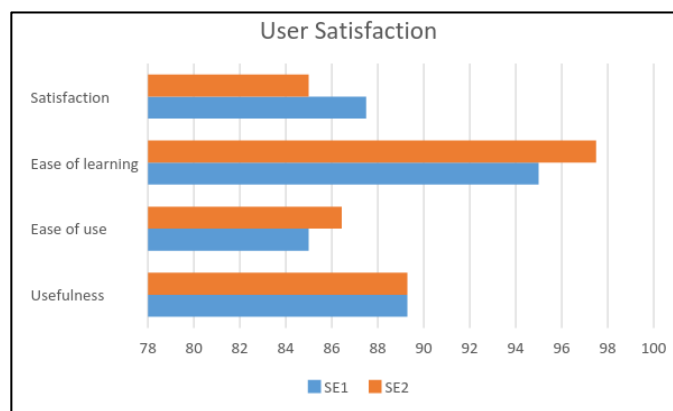| | Task Success (success/partial success/failure) | |
|---|---|---|
| | SE1 | SE2 |
| User 1 | Success | Success |
| User 2 | Success | Success |
| User 3 | Success | Partial Success |
| User 4 | Success | Success |
| User 5 | Partial Success | Success |

**Figure 6.3: Effectiveness**



**Figure 6.4: User Satisfaction Summary**

# 7      Challenges faced

- Crawling only relevant documents and discarding the irrelevant
- Creating the gold standard manually as a baseline for auto-labeling of pages

- Ensuring the learning algorithm for automatic labeling improved with new incoming pages
- Coming up the robust architecture to handle end-to-end components of the search engine
- Understanding the Elasticsearch and leveraging the correct and relevant APIs
- Creating appropriate mapping and indexing for the documents

**8      Conclusion and Future work**

The user evaluation shows that both the search engines are on-par in terms of effectiveness, efficiency and user satisfaction. The future scope involves:

- Using query logs to boost search results
- Integrate with maps
- Query suggestions
- Expand city facets beyond California

**References**

[1] https://www.elastic.co/guide/en/elasticsearch/reference/current/index.html

[2] https://www.playframework.com/documentation/2.5.x/JavaHome