



COMP8547-Advanced Computing Concepts

Project-FALL 2022

Real Estate Web Search Engine

Instructor: Dr. Mahdi Firoozjaei

Team Crawlers - Roles

Richa Singh 110093085

- Spell Checker

Sreekanth Vinodkumar 110093577

- Searching based on Bedrooms
- Recent Search History

Tanmay Verma 110089806

- Hashing
- Sorting
- Listing

Venkata Naveen Varma Vegesna 110090483

- Crawling
- HTML to Text Conversion

Venkata Rahul Nittala 110094777

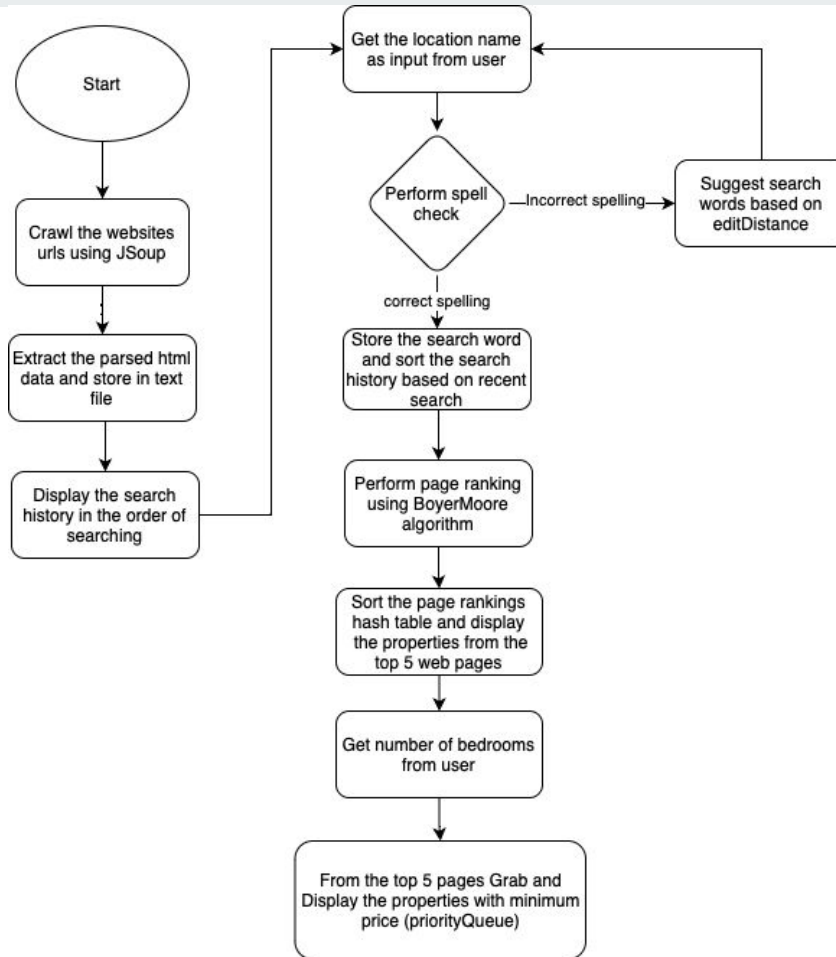
- Page Ranking

Introduction

Real Estate search engine provides the listings of properties based on the property location, number of bedrooms and price. To achieve this goal our search engine will process through several web pages from the real estate website called “<https://www.royallepagebinder.com/>” to find matches to the user's search inputs.



Workflow



Features

Crawling

Web Crawling is browsing the Web systematically to retrieve information without user intervention.

- Scrape and parse HTML from a URL.
- Crawling n number of URLs and adding the valid URLs to the HashSet and storing the html data into a “.html” file.

```
public class Crawler {  
  
    static HashSet<String> uniqueLinks = new HashSet<String>();  
    private static int max_pages = 50;  
  
    // Get links from the given url website  
    public static void crawl(String url, int page_count)  
    {  
        try {  
            if(!uniqueLinks.contains(url) && uniqueLinks.size() < max_pages) {  
                if(page_count>0) {  
                    // Add the url into the hashSet  
                    uniqueLinks.add(url);  
                    // download the html data and store as html files, from the webpage url  
                    String currentUrl = url;  
                    Document doc = Jsoup.connect(currentUrl).get();  
                    String html = doc.html();  
                    String filePath = "src/resources/htmlFiles/";  
                    // Using regular expression to format the url  
                    String fileName = currentUrl.replaceAll("[^a-zA-Z0-9_-]", "") + ".html";  
                    File tmpFile = new File(filePath);  
                    boolean fileExists = tmpFile.exists();  
                    // Check if file exists or not  
                    if(!fileExists) {  
                        // Write the formatted data into a html file.  
                        BufferedWriter out = new BufferedWriter(new FileWriter(filePath + fileName, true));  
                        out.write(url + " " + html);  
                        out.close();  
                    }  
                }  
                page_count++;  
                // Using Jsoup library to scrape and parse HTML from a URL  
                Document document= Jsoup.connect(url).get();  
                Elements linkpage= document.select("a[href^=\"https://www.royallepagebinder.com/residential-properties/page/\"]");  
                // Loop through each page link in the webpage  
                for(Element page: linkpage)  
                {  
                    crawl(page.attr("abs:href"), page_count);  
                }  
            }  
        }  
        catch(Exception e)  
        {  
            e.printStackTrace();  
        }  
    }  
}
```

HTML to Text

For each html web-page file extracted,

- Using Jsoup library, extract all property details from a web-page.
- Format the property details extracted.
- And store in a text file to make further processing easy.

```
public class Convert_to_text {  
    // Fetch html page files and scrap the data from it in required format and store it in a text file.  
    public static void convert_html_to_text()  
    {  
        System.out.println("\nFetch and convert all .html files to .txt files");  
        try {  
            String source_path = "src/resources/htmlFiles/";  
            String dest_path = "src/resources/textFiles/";  
            File dir = new File(source_path);  
            String files_list[] = dir.list();  
            String address, price, type;  
  
            // For each file, fetch property details  
            for(String filename: files_list) {  
                File currentFile = new File(source_path + filename);  
                // Using Jsoup to parse the data in the file  
                Document doc = Jsoup.parse(currentFile, "UTF-8");  
                // Fetching all the property listings in that file  
                Elements property_listings = doc.getElementsByClass("listing_quick_info");  
                // For each property, format the data as required  
                for(Element property: property_listings) {  
                    // Fetch data  
                    address = property.getElementsByClass("address").text();  
                    price = property.getElementsByClass("price").text();  
                    type = property.getElementsByClass("type").text();  
                    // Format data  
                    address = address.replace(" ", "");  
                    price = price.replace("$", "");  
                    price = price.replace(" ", "");  
                    type = type.replace(" ", "");  
                    String print_value = address+" "+price+" "+type;  
                    String fileName = filename.replace(".html", ".txt");  
                    File tmpFile = new File(dest_path);  
                    boolean fileExists = tmpFile.exists();  
                    // Check if file exists or not  
                    if(!fileExists) {  
                        // Write the formatted data into a text file.  
                        BufferedWriter out = new BufferedWriter(new FileWriter(dest_path + fileName, true));  
                        out.write(print_value + "\n");  
                        out.close();  
                    }  
                }  
            }  
            System.out.println("\nAll files converted to .txt!!\n");  
        } catch (Exception e) {  
            e.printStackTrace();  
        }  
    }  
}
```

Page Ranks

Page Rank is used to determine the rank of web page based on the comparison of the search term entered by the user in all the web pages.

- Read the downloaded web page data from the text files
- Search for the word occurrence in each file by comparing with the search term using the Boyers Moore algorithm
- Storing the count of word occurrence in the hashmap

```
public class Wordranks {
    Map<String, Integer> ranks= new HashMap<String, Integer>(); //hashmap to store the rank of word and page

    /**
     * Get the rank of the pages and the frequency of words based on the search key passed
     */
    public Map<String,Integer> getPageRank(String searchKey) {
        searchKey=searchKey.toLowerCase(); //converting the input to lower cases
        File directory=new File("src/resources/textFiles");
        File fileList[] = directory.listFiles(); //get the list of files from the directory
        /**
         * From the list of all files read each file convert to the lower case.
         */
        for(File file : fileList) {
            String line = "";
            int c=1;
            FileReader fr = null;
            try {
                fr = new FileReader(file.getPath());
            } catch (FileNotFoundException e1) {}
            e1.printStackTrace();

            Scanner scanFile=new Scanner(fr);
            while(scanFile.hasNext()) {
                line=line+" "+ scanFile.nextLine();
            }
            line=line.toLowerCase(); //convert the content to the lower case after read from file
            line=line.trim(); //remove the spaces
            int startIndex=0;
            /**
             * For each file use the Boyers Moore algorithm for searching the word count.
             */
            for(int index=0;index<=line.length();index=index+startIndex+searchKey.length()) {
                BoyerMoore bm= new BoyerMoore(searchKey);
                startIndex=bm.search(line.substring(index));
                if(startIndex+index<line.length())
                {
                    ranks.put(file.getName(), c); //store the web page name in a hash map
                    c= ranks.get(file.getName())+1; //update the frequency of the word occurrence
                }
            }
            scanFile.close();
        }
    }
}
```


Page Sorting

Page Sorting is used to sort the pages based on the ranking of the pages for the search location.

- It utilizes the Stream API of Java for sorting in descending order by using the `comparingByValue()` method and it returns a comparator that compares the values.
- Page Sorting class also gives top 5 pages out of the list which has the highest ranking.

```
1 package searchengine;
2
3 import java.util.Comparator;
4
5 public class PageSorting {
6
7     public static List<Entry<String, Integer>> SortMap(Map<String, Integer> webpages) {
8
9         List<Entry<String, Integer>> list = null;
10
11         if (!webpages.isEmpty()) {
12
13             System.out.println("\nPages after sorting based on no. of occurrences : ");
14
15             /* Using Stream API of Java 8 for sorting Hash map values.
16              * In java 8, Map.Entry class has static method comparingByValue()
17              * to help in sorting by values. This method returns a Comparator
18              * that compares Map.Entry in natural order on values.*/
19
20             Stream<Map.Entry<String, Integer>> sorted = webpages.entrySet().stream()
21                 .sorted(Map.Entry.comparingByValue(Comparator.reverseOrder()));
22
23             System.out.println("Top 5 pages with the best results are : -");
24
25             list = sorted.toList();
26
27             for (int k = 0; k < 5; k++) {
28                 System.out.println(
29                     "Page Name : " + list.get(k).getKey() + "\nNo. of Occurrences : " + list.get(k).getValue());
30             }
31
32         }
33
34         return list; // Returning the sorted list of the top 5 pages with highest word ranking
35     }
36 }
```

Top Listing

Top Listing takes the input as the Top 5 pages sorted by the page sorting feature and then reads the text files for those top 5 pages.

- It then processes that data to format, standardise and filter the data.
- Once filtered the data is stored in a Hashmap which takes Address of the listing as the Key and then Price, No. of bedroom, No. of Bathrooms and Category in a ArrayList and then passes the arraylist as the value for the corresponding key.
- This gives the data for the top listings from top 5 pages for the search location and also passes that further for searching the best deal.

```
1 package searchengine;
2
3 import java.io.BufferedWriter;
4
19
20 /* Class to get the top 5 pages by word rank and give the top listings for the Search Location*/
21
22 public class TopListing {
23
24     static HashMap<String, ArrayList<String>> ListingMap = new HashMap<>();
25
26     public static void Listing(Map<String, Integer> searchPages,String SearchLocation)
27     {
28
29         String Text_file_path = "src/resources/textFiles/";
30
31         File dir = new File(Text_file_path);
32
33         String Text_files_list[] = dir.list();
34
35         for(String filename: Text_files_list) {
36
37             if(!filename.equalsIgnoreCase(".DS_Store") && searchPages.containsKey(filename))
38             {
39
40                 try {
41                     File myObj = new File(Text_file_path+"/"+filename);
42                     Scanner myReader = new Scanner(myObj);
43                     while (myReader.hasNextLine()) {
44                         String data = myReader.nextLine();
45                         if(!data.isEmpty())
46                         {
47
48                             String[] Listings= data.split("\n");
49                             int len= Listings.length;
50                             String ListingItem="";
51                             String[] item= new String[len];
52                             String Address="";
53                             String Price="";
54                             String LastString="";
55                             String NoOfBedRooms="";
56                             String NoOfBathrooms="";
57                             String Category="";
58
59                             ArrayList<String> CombinedData=new ArrayList<>();
60                             CombinedData.clear();
61
62
63                             String[] LastStringArray= new String[3];
64
65                             /* Iterating through the text files and formatting data and storing in
66                                * Hashman and ArrayList */
67
68                             for(int i=0; i < len ; i++)
69                             {
70                                 ListingItem=Listings[i];
71                                 item=ListingItem.split(",");
72                                 Address=item[0];
```

Search History

- Save search words to database each time when user searching a word.
- Returns search words based on recent searches.

```
1 package searchengine;
2
3 import java.io.BufferedReader;
4
5 public class SearchHistory {
6     String filePath = "src/resources/searchHistory/searchHistory.txt";
7
8     /**
9      * @param word {search word}
10     */
11     public void updateSearchHistory(String word) {
12         List<String> list = new ArrayList<>();
13
14         try {
15             BufferedReader reader = new BufferedReader(new FileReader(filePath));
16             BufferedWriter writer = new BufferedWriter(new FileWriter(filePath, true));
17
18             String line = reader.readLine();
19
20             // read search words from search history file and add it to the list
21             while (line != null) {
22                 String trimmedLine = line.trim();
23                 list.add(trimmedLine);
24                 // read next line
25                 line = reader.readLine();
26             }
27             reader.close();
28
29             /**
30              * if the list contains current search word then remove the word from text file
31              * and add it to the end of the text file
32              */
33             if (list.contains(word.toLowerCase())) {
34                 clearHistory();
35                 list.forEach(item -> {
36                     try {
37                         if (!item.equals(word.toLowerCase())) {
38                             writer.write(item + "\n");
39                         }
40                     } catch (IOException e) {
41                         e.printStackTrace();
42                     }
43                 });
44                 writer.write(word.toLowerCase() + "\n");
45             } else {
46                 // else add it to the end of the text file
47                 writer.write(word.toLowerCase() + "\n");
48             }
49             writer.close();
50
51         } catch (IOException e) {
52             System.out.println("An error occurred.");
53             e.printStackTrace();
54         }
55     }
56 }
```

Best deal based on number of bedrooms

- Passing filtered data to priority queue. Filtered data contains building details in a specified location that the user has looked for.
- Generating priority queue based on price.
- Priority queue contains list of objects with details of building like address, price, number of bedrooms, and number of washrooms.
- User will be able to get plot with best price based on the number of bedrooms.

```
/**
 *
 * @param bedroomNo {number of bedrooms}
 * @return String {details of best deal}
 */
public static String bestDeal(int bedroomNo) {
    while (!buildingPriorityQueue.isEmpty()) {
        Buildings building = buildingPriorityQueue.poll();

        if (bedroomNo == building.bedrooms) {
            return "Address: " + building.address + ",\nPrice: $" + building.price + ",\nNumber of bedrooms: "
                + building.bedrooms + ",\nNumber of washrooms: " + building.washrooms;
        }
    }
    return "Not able to find a residence with given number of bedrooms";
}
```

Spell Checking

- A spell checker is a feature that is used to check the misspellings of the text.
- For performing spell check in our project we are using Edit Distance algorithm.
- Creating dictionary from the web pages.
- Using the dictionary to calculate the edit distance

```
1 package searchengine;
2
3 import java.util.Random;
4
5 public class Sequences {
6
7     public static int editDistance(String word1, String word2) {
8         int len1 = word1.length();
9         int len2 = word2.length();
10
11         // len1+1, len2+1, because finally return dp[len1][len2]
12         int[][] dp = new int[len1 + 1][len2 + 1];
13
14         for (int i = 0; i <= len1; i++) {
15             dp[i][0] = i;
16         }
17
18         for (int j = 0; j <= len2; j++) {
19             dp[0][j] = j;
20         }
21
22         // iterate though, and check last char
23         for (int i = 0; i < len1; i++) {
24             char c1 = word1.charAt(i);
25             for (int j = 0; j < len2; j++) {
26                 char c2 = word2.charAt(j);
27
28                 // if last two chars equal
29                 if (c1 == c2) {
30                     // update dp value for +1 length
31                     dp[i + 1][j + 1] = dp[i][j];
32                 } else {
33                     int replace = dp[i][j] + 1;
34                     int insert = dp[i][j + 1] + 1;
35                     int delete = dp[i + 1][j] + 1;
36
37                     int min = replace > insert ? insert : replace;
38                     min = delete > min ? min : delete;
39                     dp[i + 1][j + 1] = min;
40                 }
41             }
42         }
43
44         return dp[len1][len2];
45     }
46 } // class
```

DEMO

Github Link: https://github.com/venkata-naveen-varma-v/ACC_Project_Team_Crawlers

Thank You