NAME MEMEBR 1 : FEBRONIA ASHRAF TAWFIK                    ID : 29

NAME MEMEBR 2 : MERNA MUSTAFA EL_REFAE                    ID : 54

NAME MEMEBR 3 : NEVEEN SAMIR NAGY                         ID : 58

# LAB 3
# B TREE AND INDEXING

## ✿ problem statement :

implementing a B-tree and a simple search engine application that utilizes the B-Tree for data indexing.

### ➡ B-Tree

B-trees are balanced search trees designed to work well on disks or other direct access secondary storage devices. Unlike the red-black trees, B-tree nodes can store multiple keys and have many children. If an internal B-tree node x contains x.n keys, then x has x.n + 1 children. The keys in node x serve as dividing points separating the range of keys handled by x into x.n + 1 sub-ranges, each handled by one child of x.
Please check the reference[1] for more details about the B-Trees.

### ➡ Simple Search Engine

You will be given a set of Wikipedia documents in the XML format and you are required to implement a simple search engine that given a search query of one or multiple words you should return the matched documents and order them based on the frequency of the query words in each wiki document, please check the requirements section for more details.

---

## ✿ code design for the simple search engine application :

### 1) INDEX WEB PAGE :

Id ← id from the file
ArrayList Words ← words from this id
List<ISearchResult> list ← btree.search(every element in words)
If(list = null) → then btree.insert(word,list<ISearchResult>(id,Rank))
Else → add in this list the new element (ISearchResult (id,Rank))
END

```java
30        public void indexWebPage(String filePath) {
31            if (filePath == null || filePath == "" || !new File(filePath).exists()) {
32                throw new RuntimeErrorException(null);
33            }
34            DocumentBuilderFactory factory = DocumentBuilderFactory.newInstance();
35            DocumentBuilder builder = null;
36            try {
37                builder = factory.newDocumentBuilder();
38            } catch (ParserConfigurationException e1) {
39                // TODO Auto-generated catch block
40                e1.printStackTrace();
41            }
42            try {
43                Document document = builder.parse(filePath);
44                document.getDocumentElement().normalize();
45                // Element root = document.getDocumentElement();
46                NodeList nList = document.getElementsByTagName("doc");
47                for (int temp = 0; temp < nList.getLength(); temp++) {
48                    Node node = nList.item(temp);
49                    // System.out.println(""); //Just a separator
50                    if (node.getNodeType() == Node.ELEMENT_NODE) {
51                        // Print each employee's detail
52                        HashMap<String, Integer> map = new HashMap<String, Integer>();
53                        Element eElement = (Element) node;
54                        String id = eElement.getAttribute("id");
55                        String value = eElement.getTextContent();
56                        value = value.toLowerCase();
57                        ArrayList<String> words = new ArrayList<String>();
58                        StringTokenizer t = new StringTokenizer(value);
59                        String word = "";
60                        while (t.hasMoreTokens()) {
61                            word = t.nextToken();
62                            words.add(word.toLowerCase());
63                        }
64                        for (int j=0; j<words.size();j++) {
65                            Integer i = map.get(words.get(j));
66                            if (i == null) {
67                                map.put(words.get(j), 1);
```

```java
                        }
                        for (int j=0; j<words.size();j++) {
                            Integer i = map.get(words.get(j));
                            if (i == null) {
                                map.put(words.get(j), 1);
                            } else {
                                map.put(words.get(j), i + 1);
                            }
                        }
                        for (Entry<String, Integer> entry : map.entrySet()) {
                            String k = entry.getKey();
                            Integer v = entry.getValue();
                            List<ISearchResult> list = btree.search(k);
                            ISearchResult e = new SearchResult(id, v);
                            if (list == null) {
                                List<ISearchResult> newlist = new ArrayList<ISearchResult>();
                                newlist.add(e);
                                btree.insert(k, newlist);
                            } else {
                                list.add(e);
                            }
                        }
                    }
                }
            } catch (SAXException | IOException e) {
                // TODO Auto-generated catch block
                e.printStackTrace();
            }
        }
```

## 2) INDEX DIRECTORY :

Files[] files = Directory.listFiles()
indexWebPage(every element in files)

```java
    @Override
    public void indexDirectory(String directoryPath) {
        if (directoryPath == null || directoryPath == "" || !new File(directoryPath).exists()) {
            throw new RuntimeErrorException(null);
        }
        File folder = new File(directoryPath);
        File[] listOfFiles = folder.listFiles();
        for (File file : listOfFiles) {
            if (file.isFile()) {
                indexWebPage(directoryPath+"\\"+file.getName());
            }
        }
    }
```

## 3) DELETE WEB PAGE :

Id ← id from the file
ArrayList Words ← words from this id
List<ISearchResult> list ← btree.search(every element in words)
If(list = null) → then it not found in btree.
Else → list.remove(ISearchResult which ID = id)
END

```java
111    public void deleteWebPage(String filePath) {
112        if (filePath == null || filePath == "") {
113            throw new RuntimeErrorException(null);
114        }
115        DocumentBuilderFactory factory = DocumentBuilderFactory.newInstance();
116        DocumentBuilder builder = null;
117        try {
118            builder = factory.newDocumentBuilder();
119        } catch (ParserConfigurationException e1) {
120            // TODO Auto-generated catch block
121            e1.printStackTrace();
122        }
123        try {
124            Document document = builder.parse(filePath);
125            document.getDocumentElement().normalize();
126            // Element root = document.getDocumentElement();
127            NodeList nList = document.getElementsByTagName("doc");
128            for (int temp = 0; temp < nList.getLength(); temp++) {
129                Node node = nList.item(temp);
130                // System.out.println(""); //Just a separator
131                if (node.getNodeType() == Node.ELEMENT_NODE) {
132                    // Print each employee's detail
133                    Element eElement = (Element) node;
134                    String id = eElement.getAttribute("id");
135                    String value = eElement.getTextContent();
136                    value = value.toLowerCase();
137                    ArrayList<String> words = new ArrayList<String>();
138                    StringTokenizer t = new StringTokenizer(value);
139                    String word = "";
140                    while (t.hasMoreTokens()) {
141                        word = t.nextToken();
142                        words.add(word.toLowerCase());
143                    }
144                    for(int i=0;i<words.size();i++) {
145                        List<ISearchResult> list = btree.search(words.get(i));
146                        if(list!=null) {
147                            for(int j=0;j<list.size();j++) {
148                                if(list.get(j).getId().equals(id)) {
```

```java
129                    Node node = nList.item(temp);
130                    // System.out.println(""); //Just a separator
131                    if (node.getNodeType() == Node.ELEMENT_NODE) {
132                        // Print each employee's detail
133                        Element eElement = (Element) node;
134                        String id = eElement.getAttribute("id");
135                        String value = eElement.getTextContent();
136                        value = value.toLowerCase();
137                        ArrayList<String> words = new ArrayList<String>();
138                        StringTokenizer t = new StringTokenizer(value);
139                        String word = "";
140                        while (t.hasMoreTokens()) {
141                            word = t.nextToken();
142                            words.add(word.toLowerCase());
143                        }
144                        for(int i=0;i<words.size();i++) {
145                            List<ISearchResult> list = btree.search(words.get(i));
146                            if(list!=null) {
147                                for(int j=0;j<list.size();j++) {
148                                    if(list.get(j).getId().equals(id)) {
149                                        list.remove(j);
150                                        break;
151                                    }
152                                }
153                            }
154                        }
155                    }
156                }
157            }catch (SAXException | IOException e) {
158                // TODO Auto-generated catch block
159                e.printStackTrace();
160            }
161        }
```

## 4) SEARCH BY WORD WITH RANKING :

List<ISearchResult> list = btree.search(word)

```java
@Override
public List<ISearchResult> searchByWordWithRanking(String word) {
    if (word == null) {
        throw new RuntimeErrorException(null);
    } else if (word == "") {
        return new ArrayList<>();
    }
    return btree.search(word.toLowerCase());
}
```

## 5) SEARCH BY MULTIPLE WORD WITH RANKING :

ArrayList Words ← words from sentence

List<ISearchResult> List = btree.search(for every element in wods)

If(list = null) → then return emty list

Else → put in map word and its list

Then compare the ids in every list in map with the ids in the other lists if its equal then put in list this id and the minimum rank between them

And repeat the operation until the list is finish

```java
@Override
public List<ISearchResult> searchByMultipleWordWithRanking(String sentence) {
    ArrayList<String> words = new ArrayList<String>();
    Map<String, List<ISearchResult>> map = new HashMap<String, List<ISearchResult>>();
    SearchResult search = new SearchResult("", 0);
    if (sentence == null) {
        throw new RuntimeErrorException(null);
    } else if (sentence == "") {
        return new ArrayList<>();
    }
    StringTokenizer t = new StringTokenizer(sentence);
    String word = "";
    List<ISearchResult> mulitySearch = new ArrayList<ISearchResult>();
    int count = 0;
    while (t.hasMoreTokens()) {
        word = t.nextToken();
        words.add(word.toLowerCase());
    }
    for (int i = 0; i < words.size(); i++) {
        map.put(words.get(i), searchByWordWithRanking(words.get(i)));
        if (map.get(words.get(i)).size() == 0) {
            return new ArrayList<ISearchResult>();
        }
    }
    for (int j = 0; j < map.get(words.get(0)).size(); j++) {
        count = 1;
        search.setId(map.get(words.get(0)).get(j).getId());
        search.setRank(map.get(words.get(0)).get(j).getRank());
        for (int i = 1; i < words.size(); i++) {
            for (int k = 0; k < map.get(words.get(i)).size(); k++) {
                if (search.getId().equals(map.get(words.get(i)).get(k).getId())) {
                    count++;
                    if (search.getRank() > map.get(words.get(i)).get(k).getRank()) {
                        search.setRank(map.get(words.get(i)).get(k).getRank());
                    }
                    break;
                }
            }
        }
    }
```

✵ <u>The time complexity for the provided interfaces functions of both the B-Tree and the simple search engine :</u>

## 1) B_TREE :
### a) Insertion →
Time : O (M log N) where M is the maximum number of elements in node of B_tree & N is the number of elements in B_Tree.
### b) Search →
Time : O (M log N) where M is the maximum number of elements in node of B_tree & N is the number of elements in B_Tree.
### c) Deletion →
Time : O (M log N) where M is the maximum number of elements in node of B_tree & N is the number of elements in B_Tree.
### d) Get_Root →
Time : O(1)
### e) Get_Minimum →
Time : O(1)

## 2) SEATCH ENGINE →
### a) Index Web Page →
Time : O(I*W*M log N) where I is the number of IDs in file & W is the number of words in one documents & M log N is the time of insertion in B_Tree.
### b) Index Directory →
Time : O(F*I*W*M log N) where F is the number of files in folder & I is the number of IDs in file & W is the number of words in one documents & M log N is the time of insertion in B_Tree.
### c) Delete Web Page →
Time : O(I*W*M log N) where I is the number of IDs in file & W is the number of words in one documents & M log N is the time of deletion in B_Tree.
### d) Search by Word →
Time : O(M log N) where M log N is the time of search in B_Tree.
### e) Search by multiple word →
Time : O(L^2*W*M log N) where L is the size of list of word & W is the number of words in sentence & M log N is the time of search in B_Tree.