



**THE AMERICAN
UNIVERSITY IN CAIRO**
الجامعة الأمريكية بالقاهرة

CSCE 1101
Section 1
Dr. Howaida Ismaeel

Jana Elmidany - 900226145
Rana Taher - 900221430
Haya Shalaby - 900222031
Final Project Report

**Simple Plagiarism Detection
Utility using String Matching**

CSCE 1101-01, 02, 05 Spring 2023 Term Project Report

Department of Computer Science and Engineering, AUC

Keywords:

String-matching, Brute Force, Rabin Karp, Hamming Distance, QT, Plagiarism detection

1. Introduction

This project uses string matching algorithms Brute Force (Hamming Distance) and Rabin Karp, and their sliding window approach to identify similarity and, in turn, plagiarism. This is done by comparing a potentially plagiarized file and original files from a database.

2. Topic Definition

This report aims to examine both string-matching algorithms and compare their effectiveness. This will be done by comparing their accuracy and runtime. It will also explain the code and identify any shortcomings and limitations found in the program.

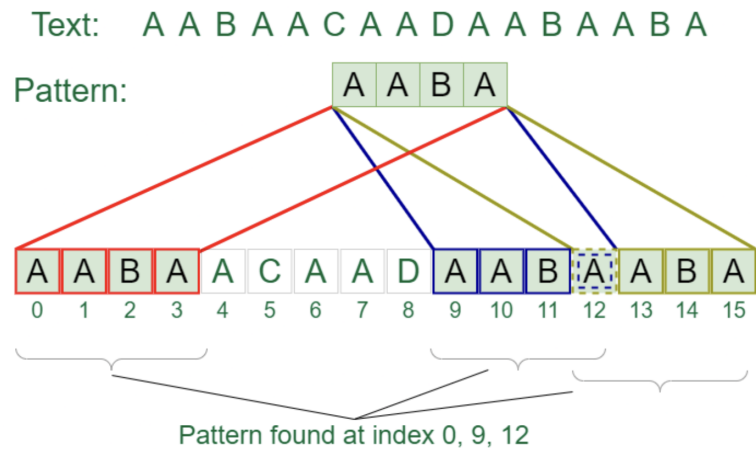
3. Methodology

Two classes were created: Document class, and PlagiarismChecker class. A PlagiarismChecker object receives the file names of the potentially plagiarized file, and the database files the user wants to check against. Creating the object calls the Document constructor which opens all these files individually and stores their lines in individual vectors. The PlagiarismChecker constructor calls all the core functions of the program which use the aforementioned algorithms to detect areas of similarity, references, quotes, and plagiarism, and increments their assigned variables accordingly.

4. Specification of Algorithms to be used

The user is given a choice between the two algorithms which are Brute Force (Hamming Distance) and Rabin Karp. Brute Force uses a function of hamming distance that takes two strings as parameters: one is the line from the user document that is checked for plagiarism and the other is a line from the database checked against. This function will count the number of

differences between the user sentence and the database sentence. Then, it will compare the differences to a threshold that we set as half of the length of the shorter line, be it the line from the user document or the line from one of the documents in the comparison database. If the differences in the characters in the lines are less than or equal to the threshold then the sentence is plagiarized, otherwise, the sentence is not plagiarized. For Rabin Karp, the sliding window mechanism is used to create a pattern, and it looks for the hash value of the pattern made from the user line in the line obtained from one of the documents from the database. We implemented the sliding window mechanisms using the help of vectors. First, we separated each line from the user document into a vector of words. Then this vector is used to create a vector of all the possible patterns that a sentence can have. Afterward, in a loop, each pattern is tested against each line in the database document. If the pattern matches any of the hash values of any segment of the lines in the database then the whole line is plagiarized, otherwise it is not.



5. Data Specifications

Three different cases need to be tested. First, we input a document to check for plagiarism that is the exact same as one of the documents in the database we are using for comparison, to show all the lines in said document are plagiarized. The second case is comparing a document to a set of documents that are completely irrelevant to it, to show that the number of lines plagiarized is supposed to be zero. The third case is when the document we are using for testing consists of some sentences from some of the documents that are used for comparison. It is supposed to show that not all of the document is plagiarized but some lines from the document are plagiarized.

6. Experimental Results

We first started by using only two files: one user file and one comparison file. We tested them being the exact same, the user file containing quotes or references, or the user file having only some copied sentences from the comparison file. When we ensured that all of these cases worked properly, we started adding and testing with more comparison files. In the terminal, we got

accurate numbers; however, when we adapted the program to Qt, the output given was inaccurate.

7. Analysis and Critique

The code we produced had several issues. The first of these is that when we are checking for the quotation in the code, it will only count the quotation if there it has one sentence. If there are two or more, it will not count them. Additionally, in-text citations affect the program's ability to detect quotes. Furthermore, spaces can cause issues when we are comparing the sentence from user and database documents. If the spacing is different, it might not be able to detect plagiarism as the lines are not aligned. Another issue we faced in the code is that if the word references or works cited were mentioned before the actual references or works cited page then it might count all the sentences after as references instead of checking them for plagiarism.

8. Conclusions

Our program is able to detect and calculate the amount of similarity, which includes references and quotes, as well as the amount of plagiarism within the user's document. After analyzing our input and code continuously, we concluded that both algorithms harbor limitations that affect their accuracy. The hamming distance method used with the brute force algorithm depends on differences which could often be misleading since plagiarism could occur on a small scale that would not be detected by this algorithm. The Rabin Karp, on the other hand, is more accurate as it actually searches for a match using hash values rather than estimating the difference between files and depending on that to reach plagiarism and similarity statistics. Apart from algorithm limitations, our own code also had a considerable amount as previously discussed. Overall, there are still many more aspects a plagiarism detection program has the potential to cover that are even more advanced such as linguistic and logical plagiarism and similarity.

Acknowledgments

Special thanks to Ahmad Abdallah for his advice and help with debugging.

References

Rabin-Karp Algorithm for Pattern Searching. (n.d.). GeeksforGeeks. Retrieved May 14, 2023, from <https://www.geeksforgeeks.org/rabin-karp-algorithm-for-pattern-searching/>

Opening and reading a file in Qt. (2017, June 23). Stack Overflow. <https://stackoverflow.com/questions/44790470/opening-and-reading-a-file-in-qt>

The Coding Train. (2019, December 18). Coding Challenge #156: Lissajous Curve Table [Video]. YouTube. https://www.youtube.com/watch?v=HS_2x8ozyAo

Appendix: Listing of all Implementation Codes

```
1 // Jana Elmidany, Haya Shalaby, Rana Taher
2 #ifndef _DOCUMENT_H
3 #define _DOCUMENT_H
4 #include <fstream>
5 #include <string>
6 #include <vector>
7 using namespace std;
8
9 class Document {
10 private:
11     string fileName;
12     vector<string> lines;
13     int numOfLines;
14
15 public:
16     Document();
17     Document(string); // constructor, opens file FN, stores lines in vector lines
18     void setFileName(string); // Sets file name
19     void setLines(string); // Sets lines
20     void addLine(string); // Adds lines
21     int numOfLines(); // Returns number of lines
22     string getFileName(); // getter for fileName
23     vector<string> getLines(); // getter for vector lines
24     //int getNumWords(); // returns no. of words in the document
25     //int countWords(string); // Counts number of words in a sentence
26 };
27
28 #endif // _DOCUMENT_H
29
30
```

```

1  #include "Document.h"
2  #include <vector>
3  #include <iostream>
4  using namespace std;
5
6  Document::Document() {
7      fileName = "";
8      numOfLines = 0;
9  }
10
11 Document::Document(string FN) {
12     setFileName(FN);
13     setLines(FN);
14 }
15
16 void Document::setFileName(string FN) { fileName = FN; }
17
18 void Document::setLines(string FN) {
19     fstream file;
20     file.open(FN);
21     string temp;
22     while (!file.eof()) {
23         getline(file, temp, '.');
24         lines.push_back(temp);
25     }
26     //cout << lines.size()-1 << endl;
27 }
28
29 void Document::addLine(string sentence) {
30     lines.push_back(sentence);
31 }
32
33
34 string Document::getFileName() { return fileName; }
35
36 vector<string> Document::getLines() { return lines; }
37
38 int Document::numOfLines()
39 {
40     double L=lines.size();
41     return L;
42 }
43
44

```

```

2  #ifndef _PLAGIARISMChecker_H
3  #define _PLAGIARISMChecker_H
4
5  #include <string>
6  #include <fstream>
7  #include <iostream>
8  #include <vector>
9  #include "Document.h"
10 using namespace std;
11
12 class PlagiarismChecker {
13 private:
14     int totalPlag;
15     int totalRef;
16     int totalSim;
17     int totalQuotes;
18     Document userDoc;
19     vector<Document> database;
20     vector<int> quoteIndex;
21     Document highlightDoc;
22     bool match ;
23     char type;
24 public:
25     PlagiarismChecker();
26     PlagiarismChecker(string userFile); //opens userFile and ogFile
27     bool isQuote(string line); //checks if similar line is a quotation & increments totalQuotes if true
28     bool isRef(string sent); //looks for reference page & adds count of all words starting there to totalRef
29     void calcPlag(); //calculates totalPlag by subtracting totalRef and totalQuotes from totalSim
30     void compDoc(); //compares userDoc against database for similarity, produces highlighted document
31     int getTotalPlag(); //getter for totalPlag
32     int getTotalRef(); //getter for totalRef
33     int getTotalSim(); //getter for totalSim
34     int getTotalQuotes(); //getter for totalQuotes
35     Document getHighlightDoc(); //getter for highlightDoc
36     void createCorpus(); //enters the files into the vector
37     void complines(vector <string> data , vector <string> user);
38     vector <string> getPat(string sentence); // gets the pattern of the sentence
39     bool bruteForce(string pat, string txt, long threshold ); //brute force & rabin karp
40     // int countWords(string sentence);
41     int hammingDistance(string str1, string str2);
42     bool rabinKarpSearch(string pattern, string text, int q);
43     vector <string> createPat(vector <string> pats);
44     void set_type(char ty);
45

```

```

1 #include <string>
2 #include <vector>
3 #include <ctype.h>
4 #include "PlagiarismChecker.h"
5 #include "Document.h"
6 using namespace std;
7
8 // d is the number of characters in the input alphabet //for stringMatching
9 #define d 256
10
11 PlagiarismChecker::PlagiarismChecker()
12 {
13     userDoc=Document();
14
15     totalPlag = 0;
16     totalRef = 0;
17     totalSim = 0;
18     totalQuotes = 0;
19     match=false;
20 }
21
22 PlagiarismChecker::PlagiarismChecker(string userFile)
23 {
24     userDoc=Document(userFile);
25     createCorpus();
26
27 }
28
29 void PlagiarismChecker::createCorpus()
30 {
31     int size;
32     string filename;
33     Document corpus;
34
35     cout<<"Enter the number of files you want to check from "<<endl;
36     cin>>size;
37
38     cout<<"Enter names of files"<<endl;
39
40     for(int i=0; i<size ; i++)
41     {
42         cin>>filename;
43
44         corpus=Document(filename+".txt");
45         database.push_back(corpus);
46     }
47 }
48
49 bool PlagiarismChecker::isQuote(string Line) {
50
51     for(int i=0 ; i<Line.length() ; i++)
52     {
53         if(Line[i] == '"' )
54         {
55             for(int j=i ; j<Line.length() ; j++)
56             {
57                 if(Line[j]=='"')
58                     return true;
59             }
60         }
61     }
62
63     return false;
64 }

```



```

66
67 bool PlagiarismChecker::isRef(string sent)
68 {
69     if(sent.substr(0,10) == "References" || sent.substr(0,11) == "Works Cited" ) //or should it be sent == "References"
70     {
71         return true;
72     }
73     return false;
74 }
75
76 void PlagiarismChecker::calcPlag()
77 {
78     totalPlag = totalSim - (totalQuotes + totalRef);
79 }
80
81 void PlagiarismChecker::compDoc() // We will count the words for references and quotes here -rana
82 {
83     vector<string> dataLines,userLines;
84     userLines=userDoc.getLines();
85     for(int i=0 ; i<database.size() ; i++)
86     {
87         dataLines=database[i].getLines();
88         complines(dataLines,userLines);
89     }
90     cout<<"The percentage of similarity is: "<<(static_cast<double>(totalQuotes+totalPlag+totalRef)/userLines.size())*1
91     cout<<"The percentage of plagiarism is: "<<(static_cast<double>(totalPlag)/userLines.size())*100<<"%<<endl;
92     cout<<"The percentage of Quotes is: "<<(static_cast<double>(totalQuotes)/userLines.size())*100<<"%<<endl;
93     cout<<"The percentage of References is: "<<(static_cast<double>(totalRef)/userLines.size())*100<<"%<<endl;
94
95 }
96
97 vector<string> PlagiarismChecker::createPat(vector<string> patterns)
98 {
99     string temp;
100     int count;
101     vector<string> fullPats;
102
103     for(int s=0 ; s<fullPats.size()+1 ; s++)
104     {
105         temp="";
106         count=0;
107
108         while(count<5 && count<patterns.size() )
109         {
110             temp += ' '+patterns[count];
111             count++;
112         }
113         if(temp!=""){
114             fullPats.push_back(temp);}
115         else{
116             break;
117         }
118
119         patterns.erase(patterns.begin());
120         if(patterns.size()<5)
121             break;
122     }
123
124     return fullPats;
125 }
126
127 void PlagiarismChecker::set_type(char ty)
128 {
129     type=ty;
130 }

```

```

131
132 void PlagiarismChecker::complines(vector <string> data , vector <string> user) // We will count the words for references
133 {
134     long threshold;
135     bool x;
136
137     vector <string> patterns;// vector of words
138     vector <string> fullPats;// vector of patterns
139
140     for(int j=0; j<user.size(); j++)
141     {
142
143
144         x=isRef(user[j]);
145
146
147         if(x==true && totalRef==0)
148         {
149             cout<<"The remaining is refreneces"<<endl;
150             for(int in=j ; in<user.size() ; in++)
151             {
152                 totalRef ++;
153                 //totalSim ++;
154             }
155             break;
156         }
157         else if(isQuote(user[j]) && totalQuotes==0)
158         {
159             totalQuotes++;
160         }
161         else
162         {
163             patterns=getPat(user[j]);
164             if(patterns.size()<5)
165             {
166                 continue;
167             }
168
169             fullPats=createPat(patterns);
170
171             bool rabin=false;
172             if(type=='r')
173             {
174                 for(int index=0 ; index<fullPats.size() ;index++){
175                     for(int k=0 ; k<data[k].size() ; k++)
176                     {
177
178                         if(rabinKarpSearch(fullPats[index], data[k] , 97))
179                         {
180                             rabin=true;
181                             totalPlag ++;
182                             break;
183                         }
184                     }
185
186                     if(rabin)
187                         break;
188                 }
189             }
190             else/

```

```

191         for(int k=0 ; k<data[k].size() ; k++)
192         {
193
194             if(data[k].length() > user[j].length())
195                 threshold=user[j].length()/2;
196             else
197                 threshold=data[j].length()/2;
198
199
200             if(bruteforce(user[j],data[k],threshold))
201             {
202                 totalPlag ++;
203                 break;
204             }
205         }
206     }
207
208 }
209
210 }
211
212
213
214 bool PlagiarismChecker::bruteforce(string usertxt, string datatxt, long threshold )
215 {
216     int distance=hammingDistance( usertxt, datatxt);
217     return (distance <= threshold);
218 }
219
220 vector <string> PlagiarismChecker::getPat(string sentence)
221 {
222     vector <string> pattern;
223     string temp;
224
225     for(int i=0 ; i<sentence.length(); )
226     {
227         temp="";
228         if(sentence[i]!=' ')
229         {
230             i++;
231             continue;
232         }
233         for(int j=i ; j<sentence.length() ; j++)
234         {
235             if (sentence[j] == ' ')
236             {
237                 break;
238             }
239             temp += sentence[j];
240         }
241
242         pattern.push_back(temp);
243         i += temp.length();
244     }
245
246     return pattern;
247
248 }

```

```

50 int PlagiarismChecker::hammingDistance(string str1, string str2)
51 {
52     int distance = 0;
53     for (int i = 0; i < str1.length(); i++) {
54         if (str1[i] != str2[i]) {
55             distance++;
56         }
57     }
58     return distance;
59 }
60
61 bool PlagiarismChecker::rabinKarpSearch(string pat, string txt , int q)
62 {
63     long M = pat.length();
64     long N = txt.length();
65     int i, j;
66     int p = 0; // hash value for pattern
67     int t = 0; // hash value for txt
68     int h = 1;
69
70     // The value of h would be "pow(d, M-1)%q"
71     for (i = 0; i < M - 1; i++)
72         h = (h * d) % q;
73
74     // Calculate the hash value of pattern and first
75     // window of text
76     for (i = 0; i < M; i++) {
77         p = (d * p + pat[i]) % q;
78         t = (d * t + txt[i]) % q;
79     }
80
81     // Slide the pattern over text one by one
82     for (i = 0; i <= N - M; i++) {
83
84         // Check the hash values of current window of text
85         // and pattern. If the hash values match then only
86         // check for characters one by one
87         if (p == t) {
88             /* Check for characters one by one */
89             for (j = 0; j < M; j++) {
90                 if (txt[i + j] != pat[j]) {
91                     break;
92                 }
93             }
94
95             // if p == t and pat[0...M-1] = txt[i, i+1,
96             // ...i+M-1]
97
98             if (j == M)
99             {
100                 cout << "Pattern found at index " << i
101                     << endl;
102                 return true;
103             }
104         }
105
106         // Calculate hash value for next window of text:
107         // Remove leading digit, add trailing digit
108         if (i < N - M) {
109             t = (d * (t - txt[i] * h) + txt[i + M]) % q;
110
111             // We might get negative value of t, converting
112             // it to positive
113             if (t < 0)
114                 t = (t + q);
115         }

```

```

305
306         // Calculate hash value for next window of text:
307         // Remove leading digit, add trailing digit
308         if (i < N - M) {
309             t = (d * (t - txt[i] * h) + txt[i + M]) % q;
310
311             // We might get negative value of t, converting
312             // it to positive
313             if (t < 0)
314                 t = (t + q);
315         }
316     }
317     return false;
318 }
319
320
321 int PlagiarismChecker::getTotalPlag()
322 {
323     return totalPlag;
324 }
325
326 int PlagiarismChecker::getTotalRef()
327 {
328     return totalRef;
329 }
330
331 int PlagiarismChecker::getTotalSim()
332 {
333     return totalSim;
334 }
335
336 int PlagiarismChecker::getTotalQuotes()
337 {
338     return totalQuotes;
339 }
340
341
342
343
344

```