

# Report: Text Extraction and Cleaning

Team Member: Ankita, Narayan, Shravan, Shadab, Suraj

## Objective:

The objective of this project is to perform text extraction and cleaning using Java programming. The main tasks involve extracting text from input sources and applying specific rules to clean the text, removing special characters, handling punctuation marks, adjusting capitalization, and eliminating unnecessary spaces and lines.

## Requirements:

Libraries and files required for the Java project:

- java.io.\*
- javax.xml.parsers.DocumentBuilderFactory
- javax.xml.parsers.DocumentBuilder
- org.w3c.dom.Document
- org.w3c.dom.NodeList
- org.w3c.dom.Node
- java.nio.file.Files
- java.io.FileWriter

## Cleaning rules for Java project:

1. If any special character other than punctuation marks is found then it will be removed completely from the text.

An array of characters to be removed fully:

```
char[] removable = {'#', '@', '$', '%', '^', '*', '~', '!', '+', '=', '_', '|', '/', '>', '<', '(', ')', '{', '}', '[', ']', '\'}
```

2. Other punctuation marks will be retained or removed as per the following conditions:

An array of punctuation marks: `char[] punctuations = {'!', ':', ';', '"', '\\", '?', '!', '!', '!'}`

- a. If a punctuation mark is between two consecutive letters or digits then it is completely removed and letters or digits are concatenated.
- b. If two or more punctuations are consecutive then the last punctuation mark will be considered and the rest are ignored.

- c. If after a character there is a punctuation mark followed by a space then it will be retained.
- 3. Numbers and letters will be retained as per the following rules:
  - a. If a digit is between two consecutive letters then it will be completely removed and letters will be concatenated.
  - b. If a letter is between two consecutive digits then it will be completely removed and digits will be concatenated.
  - c. If space is followed by digit and digit is followed by letter then remove the digit completely.
  - d. If letter is followed by digit and digit is followed by punctuation mark then remove digit completely.
- 4. Capitalization of the First letter of the sentence as per the following rules:
  - char[] delimiters = {'.', '?'};
  - a. The starting letter of the document will be capitalized.
  - b. If delimiters are followed by one space then ignore the space. (implement carefully)
  - c. If there is a delimiter before the letter then capitalize the letter.
- 5. Cleaning of unnecessary spaces and lines will be done as per the following conditions:
  - a. Consecutive two or more than two spaces will be terminated to one space.
  - b. Consecutive two or more than two lines will be terminated to one line.

### Implementation Details:

```

1  /*
2  This code is for extracting the texts from xml and text files and
3  cleaning the data as per defined rules.
4
5  Authors: Ankita, Suraj, Shadab, Shrawan, Narayan
6  Date: 4 December, 2023
7  */
8
9  // Importing necessary modules
10 import java.io.*;
11 import javax.xml.parsers.DocumentBuilderFactory;
12 import javax.xml.parsers.DocumentBuilder;
13 import org.w3c.dom.Document;
14 import org.w3c.dom.NodeList;
15 import org.w3c.dom.Node;
16 import java.nio.file.Files;
17 import java.io.FileWriter;
18
19
20 public class TextCleaner{
21     // Main method
22     public static void main(String[] args) throws IOException {
23
24         String path = "file.txt";
25         File file = new File(path);
26         FileWriter fw = new FileWriter("output.txt"); //writing to output file
27         String fileextension = extension(path);
28         String text = "";
29         if (fileextension.equals(".txt")){
30             // txt file extraction
31             text = readFileAsString(file);
32         }
33         else if (fileextension.equals(".xml")){
34             text = XMLParser(file); // Extracting text from XML file
35         }
36         if (text.equals("")){
37             System.out.println("There is no text to clean, please enter valid text.");
38         }
39         else{
40             Cleaner obj = new Cleaner();
41             String cleanedText = obj.clean(text);
42             for (int i = 0; i < cleanedText.length(); i++){
43                 fw.write(cleanedText.charAt(i));
44             }
45             fw.close();
46         }
47     }
48
49     // Identify extension
50     static String extension(String path){
51         return path.substring(path.indexOf('.'));
52     }

```

```

53
54 // txt extraction
55 public static String readFileAsString(File file) throws IOException {
56     return new String(Files.readAllBytes(file.toPath()));
57 }
58
59
60 // XML Parser Method
61 public static String XMLParser(File file){
62     try{
63         DocumentBuilderFactory dbFactory = DocumentBuilderFactory.newInstance();
64         DocumentBuilder dBuilder = dbFactory.newDocumentBuilder();
65         Document doc = dBuilder.parse(file);
66         doc.getDocumentElement().normalize();
67
68         NodeList nodeList = doc.getElementsByTagName("**"); // Get all elements
69
70         StringBuilder textContent = new StringBuilder();
71         for (int temp = 0; temp < nodeList.getLength(); temp++) {
72             Node node = nodeList.item(temp);
73             if (node.getNodeType() == Node.ELEMENT_NODE) {
74                 textContent.append(node.getTextContent()).append("\n");
75             }
76         }
77         return textContent.toString();
78     }
79     catch (Exception e) {
80         e.printStackTrace();
81     }
82     return "";
83 }
84 }
85
86
87 // Other cleaner class
88 class Cleaner{
89     char[] removable = {'#', '@', '$', '%', '^', '*', '~', '`', '+', '=', '_', '|', '/', '>', '<', '(', ')', '{', '}', '[', ']', '&'};
90     char[] punctuations = {'!', ':', ';', '?', '.', ','};
91     char[] spaces = {'\n', ' '};
92     char[] delimenaters = {'.', '?'};
93
94     // Methods to clean the data
95     public String clean(String text){
96         String t5 = text;
97         for(int i = 0; i < 3; i++){
98             String t1 = removeSpecChar(t5);
99             String t2 = updatePunct(t1);
100             String t3 = removeDigits(t2);
101             String t4 = capitalize(t3);
102             t5 = removeSpaces(t4);
103         }
104         return t5;
105     }
106 }

```

```

107 public String removeSpecChar(String text) {
108     String updatedText = "";
109     for (int i=0; i<=text.length()-1; i++){
110         char c = text.charAt(i);
111         boolean found = isCharInArray(c, removable);
112         if (found){
113             continue;
114         }
115         else{
116             updatedText += c;
117         }
118     }
119     return updatedText;
120 }
121
122
123 // Method to remove unnecessary punctuations
124 public String updatePunct(String text){
125     String updatedText = "";
126     for (int i=0; i<=text.length()-1; i++){
127         char c = text.charAt(i);
128         boolean found = isCharInArray(c, punctuations);
129         if (found){
130             if (i==0){
131                 continue;
132             }
133             else if(i==text.length()-1){
134                 updatedText += c;
135             }
136             else{
137                 char prev = text.charAt(i-1);
138                 char next = text.charAt(i+1);
139                 boolean p = isCharInArray(next, punctuations);
140                 if (Character.isLetter(prev) && Character.isLetter(next)){
141                     continue;
142                 }
143                 else if (Character.isDigit(prev) && Character.isDigit(next)){
144                     continue;
145                 }
146                 else if (Character.isDigit(prev) && Character.isLetter(next)){
147                     continue;
148                 }
149                 else if (Character.isLetter(prev) && Character.isDigit(next)){
150                     continue;
151                 }
152                 else if (p){
153                     continue;
154                 }
155                 else if (prev == ' ' && next == ' '){
156                     continue;
157                 }
158                 else if((Character.isLetter(prev) || Character.isDigit(prev)) && next == ' '){
159                     updatedText += c;

```

```

160     }
161     // in doubt unnecessary
162     else if (!(Character.isLetter(prev)) && Character.isLetter(next)){
163         if(c==',' ){
164             updatedText = updatedText + c + " ";
165         }
166         else{
167             updatedText += c;
168         }
169     }
170 }
171 else{
172     updatedText += c;
173 }
174 }
175 return updatedText;
176 }
177
178 // Combining spaces.
179 public String removeDigits(String text){
180     String updatedText = "";
181     int len = text.length();
182     for (int i=0; i<=text.length()-1; i++){
183         char c = text.charAt(i);
184         if (i == 0 && Character.isDigit(c) && Character.isLetter(text.charAt(i + 1))){
185             continue;
186         }
187         else if(i == 0){
188             updatedText +=c;
189         }
190         else if (i == len - 1 && Character.isDigit(c) && Character.isLetter(text.charAt(i - 1))){
191             continue;
192         }
193         else if(i == len -1){
194             updatedText +=c;
195         }
196         else{
197             char prev = text.charAt(i-1);
198             char next = text.charAt(i+1);
199             boolean p = isCharInArray(next, punctuations);
200             if (Character.isDigit(c) && Character.isLetter(prev) && Character.isLetter(next)){
201                 continue;
202             }
203             else if(Character.isLetter(c) && Character.isDigit(prev) && Character.isDigit(next)){
204                 continue;
205             }
206             else if (prev == ' ' && Character.isDigit(c) && Character.isLetter(next)){
207                 continue;
208             }
209             else if (next == ' ' && Character.isDigit(c) && Character.isLetter(prev)){
210                 continue;
211             }
212             else if (Character.isLetter(prev) && Character.isDigit(c) && p){

```

```

213         continue;
214     }
215     else{
216         updatedText += c;
217     }
218 }
219 }
220 return updatedText;
221 }
222
223 public String capitalize(String text){
224     String updatedText = "";
225     int len = text.length();
226     for (int i=0; i<=text.length()-1; i++){
227         char c = text.charAt(i);
228         if (i==0 && Character.isLetter(c)){
229             updatedText += Character.toUpperCase(c);
230         }else if (i == len -1 && Character.isLetter(c)){
231             updatedText = updatedText + c + ".";
232         }
233         else if (i == 0){
234             updatedText += c;
235         }
236         else{
237             char prev = text.charAt(i-1);
238             boolean p = isCharInArray(prev, delimenaters);
239             if (c == ' ' && p){
240                 continue;
241             }
242             else if ( i >= 2 && isCharInArray(text.charAt(i - 2), delimenaters) && prev == ' ' && Character.isLetter(c)){
243                 updatedText = updatedText + " " + Character.toUpperCase(c);
244             }
245             else{
246                 updatedText += c;
247             }
248         }
249     }
250     return updatedText;
251 }
252
253 public String removeSpaces(String text){
254     String updatedText = "";
255     int len = text.length();
256     for (int i=0; i<=text.length()-1; i++){
257         char c = text.charAt(i);
258         boolean p = isCharInArray(c, spaces);
259         if (i == 0 && p){
260             continue;
261         }else if(i == len - 1 && p){
262             continue;
263         }
264         else if(i == len -1){
265             updatedText +=c;

```

```

265         updatedText +=c;
266     }
267     else{
268         char next = text.charAt(i+1);
269         if (c == ' ' && next == ' '){
270             continue;
271         }else if(c == '\n' && next == '\n'){
272             continue;
273         }else{
274             updatedText += c;
275         }
276     }
277 }
278 return updatedText;
279 }
280
281 // ----- helper method-----
282 public boolean isCharInArray(char c, char[] arr){
283     boolean found = false;
284     for (char x : arr) {
285         if (x == c) {
286             found = true;
287             break;
288         }
289     }
290     return found;
291 }
292 }

```

## Result:

Input→ “hey, how are you? buddy. are you do2ne w8ith pro&j@ct. tell m#e as soon a5 s I need to make ,, .. report\*? thank you.”

Output→ “Hey, how are youbuddyyare you done with projct. Tell me as soon a s I need to make report? Thank you.”

## Conclusion:

The implemented text extraction and cleaning process successfully achieved the defined objectives. However, there were some inaccuracies observed in the cleaning process, possibly due to certain exceptional cases not being accounted for. Further refinement and handling of edge cases might be required to enhance the accuracy and robustness of the text-cleaning algorithm.



The project successfully fulfills the objectives of text extraction and cleaning in Java, but further refinement is recommended to address the observed inaccuracies.

In the future, our team would be happy to improve the cleaning part.