# COMPSCI 1XC3 - Assignment 4

Mina Al-Barak (ID: 400513160)

April 12, 2024

# Chapter 1

# Report Details

## 1.1 Program Description

The program utilizes different functions and calculations to obtain scores accurately. There were separate functions to calculate words and separate them from punctuation, such as "good,". Additionally, a separate array was designed to read emoji icons and distinguish them from words. The program combines basic programming concepts such as string length, arithmetic operations, and loops. These were crucial to read strings, distinguish words from emoji icons, and parse through all strings in a line. More complex concepts, such as dynamic memory allocation and timer, were used to efficiently store memory and hold memory blocks while processing lexicon arrays. The first function, **to_lower** converts strings into lowercase, so it can be read by the lexicon. Then, **compute_word_score** function deals with strings ending with punctuation by iterating through characters of the string and deleting the last character of the string until no punctuation marks are detected. This was an important function to integrate as, without it, the program can't process the score many words that are followed by punctuation marks and hence giving a score of 0. **is_emoticon** function checks for emoji icons in the line. It utilizes a hardcoded array of emoji icons and uses string comparison to check if any string/s in a line are emoji icons. **tokenize_string** function tokenizes sentences by breaking them down into smaller units, thereby allocating memory and adding to array. **compute_sentence_score** function takes strings and compares them to the lexicon scores to give them new scores in the array. It then calculates the average score for each sentence. Lastly, the **main** function defines the number of arguments, opens the lexicon and sentence files, calls the previous functions to parse the files, and prints the outputs.

## 1.2 Sources

I used a combination of Prof. Pasandide's lecture notes, as well as ChatGPT to help me implement the ctype.h library which I was unfamiliar with, in addition to aiding me in programming the **tokenize_string** function, and organizing the code overall.

# Chapter 2

# Appendix

## 2.1 mySA.c code

```c
#include <stdio.h> // contains printf(), fopen(), fclose(), fgets(), and sscanf().
#include <stdlib.h> // contains malloc(), realloc(), free(), and exit()
#include <string.h> // contains strlen(), strncpy(), strcmp(), strdup(), and strcspn
    ()
#include <ctype.h> // contains isspace(), ispunct(), and tolower()
#include "mySA.h"

// structure stores the words and their sentiment score
struct words parse_line(const char *line)
{ //
    struct words new_word; // struct holds words and their new score
    new_word.word = malloc(256); // memory allocation to store words using malloc()
    if (new_word.word == NULL) // if memory allocation fails
    {
        printf("Memory allocation failed.");
        exit(1); // terminate program
    }

    // sscanf() scans sentence; it passes each string called 'line' and returns the
    score based on lexicon
    sscanf(line, "%s %f", new_word.word, &new_word.score); // %s is the string, %f is
     the floating-number score
    return new_word; //return new word & corresponding score based on lexicon
}

// this function converts strings to lowercase
void to_lower(char *str)
{
    for (int i = 0; str[i]; i++) // iterate over each character in the string
    {
        str[i] = tolower(str[i]); // tolower() converts uppercase letters to
    lowercase, so
    }                               // the string can be compared to the lexicon which
    is all lowercase
}

// this function deals with strings ending with punctuation; for example "good,", "
    FUNNY!!!"
```

```
33 float compute_word_score(const char *word, struct words *lexicon, size_t lexicon_size
       )
34 {
35     char *clean_word = strdup(word); // strdup() duplicates a string and allocates
       memory for it
36     size_t len = strlen(clean_word); // strlen() finds string length
37     while (len > 0 && ispunct(clean_word[len - 1])) // ispunct() checks if a
       character is punctuation
38     {                                              // [len -1] is the last character
        (-1 due to indexing)
39         clean_word[len - 1] = '\0'; // delete punctuation character by nullifying it
40         len--; // continue shortening the string length until ispunct() can't find
       punctuation
41     }
42
43     // convert the clean word to lowercase
44     for (size_t i = 0; i < len; i++) { // iterate over each character in the string
45         clean_word[i] = tolower(clean_word[i]); // tolower() converts uppercase
       letters to lowercase
46     }
47
48     // search for clean word in lexicon
49     float score = 0.0; // initializing score
50     for (size_t i = 0; i < lexicon_size; ++i) // iterate over each word in the
       lexicon
51     {
52         if (strcmp(clean_word, lexicon[i].word) == 0) // strcmp() compares 2 strings
       to each other ie. it
53         {                                         // checks if clean word matches
        anything in the lexicon
54             score = lexicon[i].score; // score is now defined as what it is in the
       lexicon
55             break;
56         }
57     }
58
59     free(clean_word); // free memory
60     return score;
61 }
62
63 // this function checks if we're dealing with emoji icons; for example ":D" (as to
       not confuse with punctuation)
64 int is_emoticon(const char *str)
65 {
66     // made a list of emoticons to read as strings and not as punctuation
67     const char *emoticons[] = {":)", ":(", ";)", ":D", ":P", ":-)", ":-(", ";-)", ":-
       D", ":-P"};
68     int num_emoticons = sizeof(emoticons) / sizeof(emoticons[0]); // calculates
       number of emoji icons in array; good
69                                                            // for if I'd like
       to add more emojis to the array
70     // check to see if any string matches an emoji icon
71     for (int i = 0; i < num_emoticons; ++i) // iterate over each emoji icon
72     {
73         if (strcmp(str, emoticons[i]) == 0) // compare string to emoji icons to check
        if they match any emoji icon
74         {
75             return 1; // returns true if string is an emoji icon
76         }
77     }
```

```
78        return 0; // otherwise return false
79 }
80
81 // this function tokenizes string ie. breaks down sentence into smaller parts to read
       words, punctuation & spaces
82 // typically five steps: find token length, allocate memory, copy token from string,
       nullify token, add to array
83 void tokenize_string(const char *str, char **tokens, int *num_tokens)
84 {
85     int i = 0; // initialize index
86     int token_start = -1; // initialize to -1 so it can identify when a new token
       begins (see below)
87     while (str[i]) // iterate over each character in the string
88     {
89         if (!isspace(str[i])) // isspace() checks if character is a space; we're
       checking to see if it's
90         {                        // NOT a space.
91             if (token_start == -1) // checks if it's the start of a token
92             {
93                 token_start = i; // start counting index of token
94             }
95         }
96         else // otherwise, if it's NOT the start of token
97         {
98             if (token_start != -1)
99             {
100                int token_length = i - token_start; // calculate the token length
101                char *token = malloc(token_length + 1); // allocate memory for token
102                strncpy(token, &str[token_start], token_length); // strncpy() to copy
       token from the string
103                token[token_length] = '\0'; // nullify token
104                tokens[(*num_tokens)++] = token; // add token to array
105                token_start = -1; // reset token start index
106            }
107            // check if character is a punctuation and the previous character is a
       space (indicating it's
108            // the start of an emoji icon); for example: "happy :)" reads space
       before ":)"
109            if (ispunct(str[i]) && (i == 0 || isspace(str[i - 1])))
110            {
111                int emoticon_length = 1; // set emoji icon length to 1 since they'e
       considered single characters
112                                        // by array we made earlier
113                char *emoticon = malloc(emoticon_length + 1); // allocate memory for
       the emoji icon
114                strncpy(emoticon, &str[i], emoticon_length); // copy emoji icon from
       the string
115                emoticon[emoticon_length] = '\0'; // nullify eemoji icon
116                if (is_emoticon(emoticon)) // check if it's an emoticon
117                {
118                    tokens[(*num_tokens)++] = emoticon; // add to array
119                }
120                else
121                {
122                    free(emoticon); // free memory
123                }
124            }
125        }
126        i++; // iterate through all the characters
127    }
```

4

```
128     if (token_start != -1) // if there's a token remaining
129     {
130         int token_length = i - token_start; // find token length
131         char *token = malloc(token_length + 1); // allocate memory
132         strncpy(token, &str[token_start], token_length); // copy token from string
133         token[token_length] = '\0'; // nullify
134         tokens[(*num_tokens)++] = token; // add to array
135     }
136 }
137
138 // this function computes the sentiment score of a sentence
139 float compute_sentence_score(const char *sentence, struct words *lexicon, size_t
        lexicon_size)
140 {
141     float sentence_score = 0.0; // initialize sentiment score
142     int word_count = 0; // initialize word count of sentence
143
144     // tokenize sentence using last function and calculate the scores for each word
145     char *tokens[256]; // declaring array that stores the tokens
146     int num_tokens = 0; // initialize token count
147     tokenize_string(sentence, tokens, &num_tokens); // tokenize sentence
148     for (int i = 0; i < num_tokens; i++) //iterate over each token
149     {
150         float word_score = compute_word_score(tokens[i], lexicon, lexicon_size); //
        compute score for word
151         sentence_score += word_score; // update sentence score
152         word_count++; // increment word count
153         free(tokens[i]); // free allocated memory for token
154     }
155
156     if (word_count > 0) // checks if there are words in the sentence
157     {
158         sentence_score /= word_count; // average score calculator
159     }
160     return sentence_score;
161 }
162
163 // main function checks argument number, opens lexicon & sentences file, reads
        lexicon into memory, &
164 // computes sentiment score by processing sentences
165 int main(int argc, char *argv[])
166 {
167     if (argc != 3) // checks if  number of command line arguments is not 3
168     {
169         printf("Usage: %s <lexicon_file> <sentences_file>\n", argv[0]); // print
        usage message
170         return 1;
171     }
172
173     FILE *lexicon_file = fopen(argv[1], "r"); // Open lexicon file for reading; fopen
        () opens files
174     if (lexicon_file == NULL) // if file doesn't open
175     {
176         printf("Can't open file."); // print error message
177         return 1;
178     }
179
180     size_t lexicon_capacity = 100; // initialize lexicon capacity
181     size_t lexicon_size = 0; // initialize size of lexicon
```

```c
182     struct words *lexicon = malloc(lexicon_capacity * sizeof(struct words)); //
        allocate memory for lexicon
183     if (lexicon == NULL) // if memory allocation fails
184     {
185         printf("Memory allocation error."); // print error message
186         fclose(lexicon_file); // close lexicon file using fclose()
187         return 1;
188     }
189
190     char line[256]; // declare array to store lines from the file
191     while (fgets(line, sizeof(line), lexicon_file)) // fgets() reads lines from a
        file
192     {
193         if (lexicon_size >= lexicon_capacity) // if lexicon capacity is exceeded
194         {
195             lexicon_capacity *= 2; // double the capacity
196             struct words *temp = realloc(lexicon, lexicon_capacity * sizeof(struct
        words)); // reallocate memory
197             if (temp == NULL) // if reallocation fails
198             {
199                 printf("Memory allocation error"); // print error message
200                 fclose(lexicon_file); // close lexicon file
201                 free(lexicon); // free memory allocated for lexicon
202                 return 1;
203             }
204             lexicon = temp; // update lexicon timer
205         }
206         lexicon[lexicon_size++] = parse_line(line); // parse line and store word and
        score in lexicon
207     }
208     fclose(lexicon_file); // close lexicon file
209
210     FILE *sentences_file = fopen(argv[2], "r"); // open sentences file for reading
211     if (sentences_file == NULL) // if file doesn't opwn
212     {
213         printf("Can't open file."); // print error message
214         free(lexicon); // free memory allocated for lexicon
215         return 1;
216     }
217
218     // process sentences and compute scores
219     printf("string\t\t\t\t\t\t\tsample\tscore\n"); // print header
220     printf("
        ------------------------------------------------------------------------------\n")
        ; // separator
221     int sentence_number = 1; // initialize sentence number
222     while (fgets(line, sizeof(line), sentences_file)) // reading each line in the
        sentences file
223     {
224         line[strcspn(line, "\n")] = 0; // remove processed line from further
        processing using strcspn() by
225                                       // calculating initial string length &
        nullifying
226         float score = compute_sentence_score(line, lexicon, lexicon_size); // compute
         sentence score
227         printf("%-40s\t\t\t\t%.2f\n", line, score); // print sentence and score
228         sentence_number++; // increment sentence number
229     }
230
231     fclose(sentences_file); // close sentences file
```

6

```
232    free(lexicon); // free memory allocated for lexicon
233
234    return 0;
235 }
```

## 2.2 mySA.h code

```c
1  #ifndef MYSA_H
2  #define MYSA_H
3
4  struct words {
5      char *word;
6      float score;
7      float SD;
8      int SIS_array[10];
9  };
10
11 struct words parse_line(const char *line);
12
13 #endif /* MYSA_H */
```

## 2.3 Makefile

```makefile
1  CC = gcc
2  CFLAGS = -Wall -Wextra
3
4  all: mySA
5
6  mySA: mySA.c mySA.h
7
8    $(CC) $(CFLAGS) -o mySA mySA.c
9
10 clean:
11   rm -f mySA
```