## Solution to Question No. 2 Part B:

**B 2.1 Plagiarism rule and threshold:**

The rules and threshold of plagiarism in my world are as follows:

1. If there are some important words in a document that are found to be copied in the same manner then the document is plagiarized.

2. If the document is one-to-one copied then the program will output 100% plagiarized.

3. If there are some important words in a document that are found to be copied in the same manner then the document is plagiarized

4. The threshold can be defined by the user, by default the program will calculate how much of the content in the document is plagiarized, or could be plagiarized.

## B 2.2 Pseudocode for checking plagiarized content:

**bag_words:**

```
1.          Begin

2.          Create a HashMap mymap

3.          Read the data from the file and store it in char* data

4.          Declare variable tok and initilialize it with strtok(data, ".\n\t")

5.          While (tok ≠ 0)

5.1     Convert every character of tok to lower alphabet, remove the trailing and leading spaces
        and store it in to_add

5.2     If (to_add≠ '\n' && to_add ≠ '\0')

    5.2.1   If (mymap contains to_add) then increment the value of key to_add in mymap

    5.2.2   Else add to_add to hashmap and assign 1.0 as the value to the key to_add

5.3     tok = strtok(0, ".\n\t")

6.          return mymap

7.          End normalize_words:

1.          Begin

2.          Declare total_words <- 0

3.          for (j = 0; j < length of bag ; j++)

3.1             total_words = total_words + bag -> data[j].value

4.          for (j=0; j < length of bag ; j++)

4.1             set the value for key bag -> data[j].key as current value for the key
        / total_words

5.          End
```

### compute_tf_idf:

1. Begin
2. for (i = 0 ; i < length of bag ; i++)

      1.1          curr_word = bag - > data[i].key

      1.2          tf = bag -> data[i].value

      1.3          number_of_documents_with_word = 0

      1.4          for (j = 0 ; j < length of bags ; j++)

  1.4.1          curr_bag = bags -> data[j]

  1.4.2        if    (curr_bag        contains curr_word)       then number_of_documents_with_word++

      1.5          idf = 1 + log (length of bags / number_of_documents_with_word)

      1.6          set the bag with key curr_word as tf * idf

3. End
4. **compute_cosine_similarity:**
1.        Begin
2.        modQuery = 0
3.        for (i = 0; i < length of query ; i++)

    3.1     modQuery = modQuery + (query -> data[i].value)$^2$

4.        modQuery = $\sqrt{modQuery}$
5.        modDocument = 0
6.        for (i = 0; i < length of document ; i++)

    6.1     modDocument = modDocument + (document -> data[i].value)$^2$

7.        modDocument = $\sqrt{modDocument}$
8.        querytimesdocument = 0
9.        for (i = 0; i < length of document ; i++)

    9.1     if (document)

10.      End

### compare_query_with_docum ents:

1. Begin
2. Declare Vecor bags, HashMap query_bag
3. Read all the document files and use bag_words to bag the words and store it in a HashMap and add it to the Vector bags
4. Bag the Query using bag_words and store it in query_bag
5. Normalize the bags using normalize_bags
6. Normalize the Query using normalize_bags
7. Compute TF*IDF using compute_tf_idf for bags
8. Compute TF*IDF using compute_tf_idf fro query

9. End

B 2.3 C Program:

```c
#include <stdio.h>
#include <string.h>
#include <stdbool.h>
#include <dirent.h>
#include <math.h>
#include "hash_map.h"
#include "pcheck.h"
#include "debug_helper.h"
#include "file_reader.h"
#include "input_helper.h"
#include "vector.h"

int main(int argc, char** argv) {

    char* file_name = argv[0];
    char* directory = "./data";
    char* query_file;
    float threshold = 80.0f;;

    switch(argc) {
        case 2:
            query_file = argv[1];
            break;
        case 3:
            query_file = argv[1];
            threshold = atoi(argv[2]);
            break;
        case 4:
            directory = argv[1];
            query_file = argv[2];
            threshold = atoi(argv[3]);
```

```
32                    break;
33               default:
34                   printf( "\nUSAGE: %s <directory> <file> <threshold>"
35                               "\ndefault directory = ./data"
36                               "\ndefault threshold = 80.0%%"
37                               "\nExample: %s data query 92.2\n", file_name, file_name)
38                   break;
39          }
40
41          if (argc <= 1 || argc > 4)
42              return -1;
43
44          struct dirent* in_file;
45          FILE* entry_file;
46          DIR *FD;
47
48          if ((FD = opendir(directory)) == NULL) {
49              printf("\nError Opening %s Directory\n", directory);
50              return -1;
51          }
52
53          Vector* file_list = newMinimalVector();
54
55          while (in_file = readdir(FD)) {
56              if (!strcmp (in_file->d_name, "."))
57                  continue;
58              if (!strcmp (in_file->d_name, ".."))
59                  continue;
60
```

```
60
61          //ds(in_file -> d_name);
62          char* path = new_string(directory);
63          strcat(path, "/");
64          strcat(path, in_file -> d_name);
65          file_list -> add(file_list, path);
66      }
67      //return 0;
68      Vector* bags;
69      HashMap* query_bag;
70
71      bags = newMinimalVector();
72      query_bag = newHashMap();
73
74      /* Bag the Worlds */
75      for (int i = 0 ; i < file_list -> length ; i++)
76          bags -> add(bags, bag_words(file_list -> data[i]));
77
78      /* bag the query */
79      query_bag = bag_words(query_file);
80
81      /* Normalize the Bags */
82      for (int i = 0 ; i < bags -> length ; i++)
83          normalize_words(bags -> data[i]);
84
85      /* Normalize the Query */
86      normalize_words(query_bag);
87
```

```
87
88          /* Compute tf*idf */
89          /* Make sure that you have already computed TF and also normalized it */
90          for (int i = 0 ; i < bags -> length ; i++)
91              compute_tf_idf(bags -> data[i], bags);
92
93          compute_tf_idf(query_bag, bags);
94
95          printf("\nTF-IDF COSINE SIMILARITY PLAGIARIZATION CHECKER ------------------\n\n");
96
97          /* find the cosine similarity */
98          int count = 0;
99          printf("\nQUERY --------------------------------------: \n");
100         printf("%s\n\n", read_string_file(query_file));
101         for (int i = 0 ; i < bags -> length ; i++) {
102             double similarity = compute_cosine_similarity(query_bag, bags -> data[i]);
103             if (similarity * 100 > threshold) {
104                 count++;
105                 printf("\nFOUND PLAGIARIZATION with %.10f%% Cosine similarity\n", similarity*100);
106                 printf("\nORIGINAL FILE -----------------------------------------: \n");
107                 printf("%s\n",read_string_file(file_list -> data [i]));
108             }
109         }
110
111         printf("\nSearched %d documents, found %d matches with minimum threshold of %.2f%%\n\n",
112             bags -> length, count, threshold);
113
114         return 0;
115 }
```

**pcheck.c**

```c
1    #include <stdio.h>
2    #include <string.h>
3    #include <stdbool.h>
4    #include <math.h>
5    #include "hash_map.h"
6    #include "pcheck.h"
7    #include "debug_helper.h"
8    #include "file_reader.h"
9    #include "input_helper.h"
10   #include "vector.h"
11
12   const char *delim = ".\n\t ";
13
14   HashMap *bag_words(char *file) {
15       HashMap *mymap = newHashMap();
16       char *data = read_string_file(file);
17       char *tok;
18       tok = strtok(data, delim);
19       while (tok != 0) {
20           char *to_add = strlwr(new_string(sanitize_string(tok)));
21           if (*to_add != '\n' && *to_add != '\0' && strlen(to_add) > 0) {
22               if (mymap->contains(mymap, to_add))
23                   mymap->set(mymap, to_add, new_double((*(double *)mymap->get(mymap, to_add)) + 1.0));
24               else
25                   mymap->add(mymap, to_add, new_double(1.0));
26           }
27           tok = strtok(0, delim);
28       }
29       return mymap;
30   }
```

```
32      /* found normalized TF */
33      void normalize_words(HashMap *bag) {
34          double total_words = 0;
35          for (int j = 0; j < bag->length; j++)
36              total_words += *(double *)bag->data[j].value;
37          for (int j = 0; j < bag->length; j++)
38              bag->set(bag, bag->data[j].key, new_double((*(double *)bag->data[j].value) / total_w
39      }
40
41      /* TF*IDF
42      ** IDF = 1 + log(Total number of documents / Number of documents with that word)*/
43      void compute_tf_idf(HashMap *bag, Vector *bags) {
44          for (int i = 0; i < bag->length; i++) {
45              char *curr_word = bag->data[i].key;
46              double tf = *(double *)bag->data[i].value;
47              int number_of_documents_with_word = 0;
48              for (int j = 0; j < bags->length; j++) {
49                  HashMap *curr_bag = (HashMap *)(bags->data[j]);
50                  if (curr_bag->contains(curr_bag, curr_word))
51                      number_of_documents_with_word++;
52              }
53              double idf = 1.0 + log(((double)bags->length) / number_of_documents_with_word);
54              bag->set(bag, curr_word, new_double(tf * idf));
55          }
56      }
57
```

```c
57     /* cos(theta) = A.B / (||A||*||B||)
58     ** here ||A|| = (TF*IDF for word 1)^2 + (TF*IDF for word 2)^2 and so on,
59     ** same is applied for ||B||*/
60     double compute_cosine_similarity(HashMap *query, HashMap *document) {
61         /* computer ||query|| */
62         double modQuery = 0.0;
63         for (int i = 0; i < query->length; i++)
64             modQuery += pow(*(double *)(query->data[i].value), 2);
65         modQuery = sqrt(modQuery);
66         //dd(modQuery);
67         /* compute ||document||, only take the words that are present in the query */
68         double modDocument = 0.0;
69         for (int i = 0; i < document->length; i++)
70             if (query->contains(query, document->data[i].key))
71                 modDocument += pow(*(double *)(document->data[i].value), 2);
72         modDocument = sqrt(modDocument);
73         //dd(modDocument);
74
75         /* compute query.document = q1*d1+q1*d2 . . . .qn*dn */
76         double querytimesdocument = 0.0;
77         for (int i = 0; i < query->length; i++)
78             if (document->contains(document, query->data[i].key))
79                 querytimesdocument += *(double *)(query->data[i].value) * *(double *)(document->get(document, query->data[i].key));
80         //dd(querytimesdocument);
81         double cosine_similarity = 0.0;
82         cosine_similarity = querytimesdocument / (modQuery * modDocument);
83         //dd(cosine_similarity);
84         return cosine_similarity;
85     }
```

Output

```
TF-IDF COSINE SIMILARITY PLAGIARIZATION CHECKER ------------------


QUERY ----------------------------------------:
The domestic cat (Felis silvestris catus or Felis catus) is a small, typically furry, carnivorous mammal.
They are often called house cats when kept as indoor pets or simply cats when there is no need to distinguish them from other felids and felines.


FOUND PLAGIARIZATION with 93.2975787527% Cosine similarity

ORIGINAL FILE ----------------------------------------:
The domestic cat (Felis silvestris catus or Felis catus) is a small, typically furry, carnivorous mammal.
They are often called house cats when kept as indoor pets or simply cats when there is no need to distinguish them from other felids and felines.
They are often valued by humans for companionship and for their ability to hunt vermin. There are more than seventy cat breeds recognized by various cat registries.

Searched 4 documents, found 1 matches with minimum threshold of 90.00%
```

## **Explanation:**

All the Source Documents are stored in a folder named data, which are text files, the query is also stored in a file named query that contains the potentially plagiarized document. These files are read, they are sanitized, i.e. the leading spaces, and trailing spaces are removes, they are split into words, bagged, and then each of the documents are converted to n-dimensional vector, the query is also converted to ndimensional vector, these vectors are then normalized, and the TF*IDF is calculated for each of them, the advantage of using TF*IDF is that common words across the documents are weighted down, i.e. they are assigned a lower score, and the unique words are weighted more, using logarithmic, now that we have a vector for each of the document and also the query, it can simply be compared using Dot Product, which is

$$A.B = ||A|| \times ||B|| \times \cos\theta$$

$$\cos\theta = \frac{A.B}{||A|| \times ||B||}$$

cos $\theta$ will output a value from $[-1, 1]$, but for our implementation it won't output a value less than 0, hence the output of Cosine Similarity will be from $[0, 1]$, this can be mapped from $[0, 100]$ to get a percentage of potential plagiarisation.

It then checks which documents have a minimum plagiarisation which was passed as a command line argument to the program, and those plagiarized documents are displayed along with the query.

**Bibliography** _____

1. Francis Padillo (2016) A Data Structure to Speed-Up Machine Learning Algorithms on Massive Datasets.
2. Supreethi, K.P., Prasad, E.V.: Web Document Clustering Technique Using Case Grammar Structure. In: IEEE International Conference on Computational Intelligence and Multimedia Applications 2007 (2007).
3. US Patent US8375021B2, https://patents.google.com/patent/US8375021B2/en.
4. For further source codes regarding plagiarism question refer https://github.com/Deepakr120/plagarism-detector-using-c