

Proiectarea Algoritmilor

Tema de casa

Distanța între două siruri

Chiriac Adrian-Florin

Mai 2024

1 Introducere

1.1 Enunțul problemei

We have the task of developing an algorithm for an advanced code editor that automatically corrects syntax errors in programming languages. It is assumed that we receive a clear specification of the valid syntax of the programming language in the form of a "rule" and a code fragment that contains syntax errors, i.e., does not conform to that rule.

Our objective is to build an algorithm that determines the minimum number of operations required to transform the code fragment

into one that complies with the given rule. These operations may include character substitutions, insertions, or deletions.

Let's take a concrete example to illustrate the problem: let's assume we have the following syntax rule for function declarations in

the programming language: "Every function must start with the keyword "func", followed by the function name enclosed in parentheses." An example of a valid function declaration would be "func(myFunction)". Here's how the situation looks: Given code fragment: "fnuc(myFuncion"

Our objective is to find the minimum number of operations required to correct the code fragment so that it matches the pattern

defined by the rule. These operations may include, for example, reversing the characters "n" and "u" to obtain "func", then inserting

the missing characters "t" and ")", so that we obtain "func(myFunc)" according to the given rule.

1.2 Descrierea problemei

We're given two strings and we have to determine the minimum number of operations that it takes to transform the second string into the rule string (valid one). The operations we're able to make are:

- Substitute a character
- Delete a character
- Insert a character
- Reverse two characters

2 Algoritmi

2.1 Pseudocod

```
DECLARE rule AS STRING
DECLARE code_to_change AS STRING
DECLARE n AS INTEGER
DECLARE m AS INTEGER
```

```

DECLARE dp[105][105] AS INTEGER
FUNCTION MinOfTwo(integer a, integer b)
IF a>b THEN:
RETURN a
else
RETURN b
END IF
END FUNCTION

FUNCTION MinimumNoOperations( string rule , string code_to_change, integer dp[105][105], integer n, integer m)
FOR i = 0 TO n+1 DO
dp[i][0]=i
END FOR
FOR j = 0 TO m+1 DO
dp[0][j]=j
END FOR
FOR i = 0 TO n DO
FOR j=0 TO m DO
IF rule[i]=code_to_change[j] THEN:
dp[i+1][j+1]=dp[i][j]
ELSE
integer insert_cost=dp[i+1][j]+1
integer delete_cost=dp[i][j+1]+1
integer replace_cost=dp[i][j]+1
IF rule[i-1]==code_to_change[j] and rule[i]==code_to_change[j-1] THEN:
swap_cost=dp[i-1][j-1]+1
replace_cost=MinOfTwo(replace_cost,swap_cost)
END IF
dp[i+1][j+1]=MinOfTwo(MinOfTwo(replace_cost,insert_cost),delete_cost)
END IF
END FOR
END FOR
RETURN dp[n][m]
END FUNCTION

READ rule
READ code_to_change

FOR i = 0 TO 100 DO
FOR j=0 TO 100 DO
dp[i][j]=0
END FOR
END FOR

n = length(rule)
m = length(code_to_change)
integer nr = MinimumNoOperations(rule,code_to_change,dp,n,m)
PRINT "Minimum number of operations required to transform the second string to a valid one is equal to:"
PRINT nr

```

3 Descrierea aplicatiei

3.1 Utilizare

The application functions in 2 working regimes:

- User Input
- Between 2 and 5 randomly chosen words from a list

1) For the first regime, the user is prompted with the question "Input the rule string:" and after that "Input the fragment code:". (For example, the rule string is "func(myFunction)" and the fragment code is "fnuc(myFunction)"). The program is going to output "3", because there are 3 steps to accomplish this which are explained in the document:

- Reversing the characters "n" and "u"
- Adding the letter "t"
- Adding the character ")"

Therefore there are 3 operations.

2) In the second regime, the program takes its data from a .txt file using file pointer. In the .txt file there are 1000 randomly generated words. The program then chooses a random number between 2 and 5 (this will be the number of words in the string). The individual words are then randomly chosen with an "index=rand()MOD1000" and concatenated to the rule string. The code fragment is then generated by a "similar string generator" that we will explain in a different section to make sure all the details are clear.

2*) SimilarStringGenerator Function:

Explanation of the similar_string_generator Function

Summary:

The function creates a string similar to the input string but with a random number of errors (between 1 and the length of the string). These errors can be character swaps, deletions, or replacements. This is useful for testing automatic syntax error correction algorithms by simulating different types of errors that can occur in a character string.

Step-by-step Explanation:

1. `size_t length = strlen(source);`
Calculate the length of the source string.

2. Allocate memory for the similar string:

```
1 char *similarString = malloc(length + 1);
2 if (!similarString) {
3     return NULL;
4 }
```

3. Copy the source string into the result string:

```
1 strcpy(similarString, source);
```

4. Generate the number of errors:

```
1 int num_errors = (rand() % length) + 1;
```

5. Introduce the errors:

```
1 for (int i = 0; i < num_errors; i++) {
2     int error_type = rand() % 3;
3     int pos = rand() % length;
4     if (error_type == 0 && pos < length - 1) {
5         char temp = similarString[pos];
6         similarString[pos] = similarString[pos + 1];
```

```

7         similarString[pos + 1] = temp;
8     } else if (error_type == 1 && length > 1) {
9         memmove(&similarString[pos], &similarString[pos + 1], length
10             - pos);
11         similarString[length - 1] = '\0';
12         length--;
13     } else if (error_type == 2) {
14         similarString[pos] = 'a' + (rand() % 26);
15     }
16 }

```

Iterate `num_errors` times to introduce errors into the similar string:

- **Swap:** If `error_type` is 0 and `pos` is less than `length-1`, swap two adjacent characters.
- **Deletion:** If `error_type` is 1 and the string length is greater than 1, delete the character at position `pos` and shift the characters to the right of it one position to the left, reducing the string length.
- **Replacement:** If `error_type` is 2, replace the character at position `pos` with a random character between 'a' and 'z'.

6. Return the modified string:

```

1     return similarString;

```

4 Rezultate

```

Edit Distance w/ swap:

Enter your choice
1 for user input
2 for randomly chosen words
3 for exiting the program

1

Please type the rule string:      func(myFunction)
Please type the user input string: fnuc(myFuncion
Minimum number of operations required: 3

```

Figure 1: The first example

```

2
The randomly generated word is: mother_industry_establish_keep_bed_hospital_hope_assume_task_dinner_policy_think_side_step
_maybe_hang_party_head_realize_view

The randomly generated code fragment (but still similar to rule) is: mother_industr_yestbaliss_keep_bed_hospital_hope_assm
e_task_dinner_policy_thgnk_side_stepqmayeb_hang_partyhead_realize_view

Minimum number of operations required: 125

```

Figure 2: The second example

```

2
The randomly generated word is: he_employee_pain_another_however_certainly_there_Mrs_save_out_example_organization_student
_option_raise_total_view_offer_risk_action_create_right_ready_Mrs_understand_last_possible_human_oil_why_range_board_futur
e_trade_hotel_section_into_perform_long_administration_fire_image_office_science_hour_finish_good_what_movie_eye_record_ri
se_act_poor_up_table_free_sure_son_society_painting_available_eye_explain_lead_hear_election_want_manage_involve_listen_fo
cus_its_pretty_find_me_visit_final_any_determine_teacher_voice_task_than_address_should_surface_lawyer_PM_pressure_reduce_
explain_defense_commercial_TV_worry

The randomly generated code fragment (but still similar to rule) is: he_employe_pain_another_however_crtanly_there_Mrusaveo
ut_example_ognaization_studet_optlon_raise_totla_view_ofer_irs_k_atcioubcreate_right_raedy_Mrl_understand_lasn_poshibte_uma
_loi_why_range_obyr_g_fuutre_trade_hoetl_section_into_perform_lo_nadministratio_dire_image_fofice_sciencebhour_finish_good_h
watvmovi_ey_erecordsriseiact_poor_kp_table_frees_ure_sno_soiety_apintni_available_eye_expkain_nead_hearelectio_nwant_mawa
geinvolve_listen_focus_its_lretty_find_mi_visitsfinal_any_determine_techer_pocce_task_han_aydressvshould_surace_alwuemfPM_p
ressure_reducanxpain_defenhe_commercial_TV_worry

Minimum number of operations required: 613

```

Figure 3: The third example

```

3

Process returned 0 (0x0)    execution time : 95.151 s
Press any key to continue.

```

Figure 4: The 4th example (Making sure the program closes properly)

5 Algoritmul lui Levenshtein Personalizat

We will use the Levenshtein algorithm from the famous "Edit Distance" problem, where there are only three operations (the inversion operation does not exist). We will use a dp matrix (an idea based on dynamic programming) that works as follows:

- If two characters are the same, we do nothing.
- If the two characters are different, we have four operations: deletion, insertion, substitution, inversion.

Imagine that the rows represent the first string, and the columns represent the second string. If you want to delete a character, you go back to the previous column ($dp[i][j-1]$). If you want to add a letter, you go back to the previous row ($dp[i-1][j]$). If you want to substitute the last character in both strings, you go back to $dp[i-1][j-1]$. If you want to perform an inversion, you need to go to $dp[i-2][j-2]$.

In the end, to reach $dp[i][j]$, you need to perform the minimum of these four operations and add one.

Example: Customized Levenshtein Algorithm (with inversion operation)

Target string: func(myFunction)

Code fragment: fnuc(myFunction)

To transform the code fragment into the target string, the following operations are needed:

1. Invert "n" and "u".
2. Add the letter "t".
3. Add the character ")".

So, in total, there are 3 operations.

Let n be the length of the target string, which is 16, and m be the length of the code fragment, which is 14. Our solution is $dp[n][m]$, which is $dp[16][14]$, and its value is 3. This verifies the algorithm.

6 Program C

The program is composed by 3 files:

- The main file
- The file that generates the strings
- The header file which calls the functions

6.1 Main.c

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <time.h>
4  #include <unistd.h>
5  #include "choosewords.h"
6  #define DIM 200000
7  #include <string.h>
8  char rule[DIM], input[DIM];
9  char *fragment_cod;
10 int dp[DIM][DIM], distance;
11
12 int Min(int a, int b)
13 {
14     return a>b ? b : a;
15 }
16
17 int MinimumNoOperations(char rule[DIM], char code_fragment[DIM], int dp[DIM][DIM], int n, int m)
18 {
19     for(int i=0; i<=n; i++)
```

```

20 {
21     dp[i][0]=i;
22 }
23 for(int j=0; j<=m; j++)
24 {
25     dp[0][j]=j;
26 }
27
28 for(int i=0; i<n; i++)
29 {
30     for(int j=0; j<m; j++)
31     {
32         if(rule[i]==code_fragment[j])
33             dp[i+1][j+1]=dp[i][j];
34         else
35         {
36             int insert_cost=dp[i+1][j]+1;
37             int delete_cost=dp[i][j+1]+1;
38             int replace_cost=dp[i][j]+1;
39
40             if(rule[i-1]==code_fragment[j] && rule[i]==code_fragment[j
41                 -1])
42             {
43                 int swap_cost=dp[i-1][j-1]+1; // bonus operation
44                 compared to the original "Edit Distance" problem
45                 replace_cost=Min(replace_cost,swap_cost);
46             }
47
48             dp[i+1][j+1]=Min(Min(insert_cost,delete_cost), replace_cost)
49             ;
50         }
51     }
52 }
53 return dp[n][m];
54 }
55
56 int main ()
57 {
58
59     printf("Edit Distance w/ swap:\n");
60
61     while(1)
62     {
63         printf("\nEnter your choice\n");
64         printf("1 for user input\n");
65         printf("2 for randomly chosen words\n");
66         printf("3 for exiting the program\n\n");
67
68         int choice;
69         scanf("%d",&choice);
70
71         switch (choice)
72         {
73             case 1:
74                 printf("\n");
75                 getchar();
76                 printf("Please type the rule string:      ");
77                 fgets(rule, DIM, stdin);
78                 rule[strlen(rule)-1]=0;
79
80                 printf("Please type the user input string:  ");

```

```

78         fgets(input, DIM, stdin);
79         input[strlen(input)-1]=0;
80
81         int len_input = strlen(input);
82         int len_rule = strlen(rule);
83
84         distance = MinimumNoOperations(rule, input, dp, len_rule,
            len_input);
85
86         printf("Minimum number of operations required: %d\n", distance);
87         memset(input,0,DIM);
88         memset(rule,0,DIM);
89
90         break;
91
92     case 2:
93         getchar();
94         printf("\n");
95
96         char **words= allocate2DArray(1000, 50);
97         if(words==NULL)
98         {
99             exit(4);
100        }
101
102        choose_random_word(words, rule);
103        printf("The randomly generated word is: ");
104        printf("%s",rule);
105        printf("\n\n");
106        int n=strlen(rule);
107        fragment_cod=similar_string_generator(rule);
108        printf("The randomly generated code fragment (but still similar
            to rule) is: ");
109        printf("%s",fragment_cod);
110        printf("\n\n");
111        int m=strlen(fragment_cod);
112
113        distance = MinimumNoOperations(rule, input, dp, n, m);
114        printf("Minimum number of operations required: %d\n", distance);
115
116        free2DArray(words, 1000);
117        memset(input,0,DIM);
118        memset(rule,0,DIM);
119
120        break;
121
122    case 3:
123        return 0;
124        break;
125
126    default:
127        printf("\nInvalid choice!\n");
128        break;
129    }
130
131    }
132    return 0;
133 }

```


6.2 Header

```
1 #ifndef CHOOSEWORDS_H
2 #define CHOOSEWORDS_H
3
4 void choose_random_word(char **words, char rule[]);
5 char *similar_string_generator(char *source);
6 char** allocate2DArray(int rows, int cols);
7 void free2DArray(char **array, int rows);
8 #endif // CHOOSEWORDS_H
```

6.3 Generator File

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <time.h>
4 #include <string.h>
5 #include <unistd.h>
6 #include "choosewords.h"
7 #define WORD_COUNT 10000
8
9 void choose_random_word(char **words, char rule[])
10 {
11
12     FILE *samplewords = fopen("wordlist.txt","r");
13     if( samplewords == NULL)
14     {
15         printf("the file could not be opened!");
16         exit(1);
17     }
18
19     int i=0;
20     while(fgets(words[i],40,samplewords)!= NULL)
21     {
22         words[i][strlen(words[i])-1]=0;
23         i++;
24     }
25
26     int howmanywords=i;
27
28     /*
29     printf("iata cateva cuvinte!\n");
30     for(i=0; i<howmanywords; i++)
31     {
32         printf("%s\n",words[i]);
33     }
34     */
35
36     srand(time(NULL)*getpid());
37
38     //0 sa aleg intre 2 si 5 cuvinte
39     int sentencelength=rand()%100;
40
41     for(int i=0; i<sentencelength; i++)
42     {
43         int random_index=rand()%1000;
44         strcat(rule, words[random_index]);
45         strcat(rule, "_");
46     }
47     rule[strlen(rule)-1]=0;//eliminates the last '_'
48 }
```

```

49 }
50
51 char *similar_string_generator(char *source)
52 {
53     size_t length = strlen(source);
54
55     char *similarString = malloc(length + 1);
56     if (!similarString)
57     {
58         exit(1);
59     }
60     strcpy(similarString, source);
61
62     int num_errors = (rand() % length) + 1;
63
64     for (int i = 0; i < num_errors; i++)
65     {
66         int error_type = rand() % 3;
67         int pos = rand() % length;
68
69         if (error_type == 0 && pos < length - 1)
70         {
71             char temp = similarString[pos];
72             similarString[pos] = similarString[pos + 1];
73             similarString[pos + 1] = temp;
74         }
75         else if (error_type == 1 && length > 1)
76         {
77             memmove(&similarString[pos], &similarString[pos + 1], length -
78                 pos);
79             similarString[length - 1] = '\0';
80             length--;
81         }
82         else if (error_type == 2)
83         {
84             similarString[pos] = 'a' + (rand() % 26);
85         }
86     }
87     return similarString;
88 }
89
90
91 char** allocate2DArray(int rows, int cols)
92 {
93     char **array;
94     int i;
95
96     // Allocate memory for the array of pointers
97     array = (char**) calloc(rows, sizeof(char*));
98
99     if (array == NULL)
100     {
101         printf("Memory allocation failed.\n");
102         exit(2);
103     }
104
105     // Allocate memory for each row
106     for (i = 0; i < rows; i++)
107     {
108         array[i] = (char*) calloc(cols, sizeof(char));

```

```

109         if (array[i] == NULL)
110         {
111             printf("Memory allocation failed.\n");
112
113             // Free previously allocated rows
114             while (--i >= 0)
115                 free(array[i]);
116             free(array);
117             exit(3);
118         }
119     }
120 }
121
122 return array;
123 }
124
125
126 void free2DArray(char **array, int rows)
127 {
128     for (int i = 0; i < rows; i++) {
129         free(array[i]);
130     }
131     free(array);
132 }

```

7 Program Python

The program is based on 1 file.

```

1
2 def MinOfTwo(a,b):
3     if a<b:
4         return a
5     else:
6         return b
7 def MinimumNoOperations(rule,code_fragment):
8     n=len(rule)
9     m=len(code_fragment)
10    dp = [[0] * (m+1) for _ in range(n+1)]
11    for i in range(n+1):
12        dp[i][0]=i
13    for j in range (m+1):
14        dp[0][j]=j
15    for i in range (n):
16        for j in range(m):
17            if rule[i]==code_fragment[j]:
18                dp[i+1][j+1]=dp[i][j]
19            else:
20                insert_cost=dp[i+1][j]+1
21                delete_cost=dp[i][j+1]+1
22                replace_cost=dp[i][j]+1
23                if rule[i-1]==code_fragment[j] and rule[i]==code_fragment[j
24                    -1]:
25                    swap_cost=dp[i-1][j-1]+1
26                    replace_cost=MinOfTwo(replace_cost,swap_cost)
27                dp[i+1][j+1]=MinOfTwo(MinOfTwo(insert_cost,delete_cost),
28                    replace_cost)
29    return dp[n][m]
30 rule=input("Type the rule that will define the correct syntax: ")
31 code_fragment=input("Introduce the code fragment : ")

```

```
30 nr=MinimumNoOperations(rule,code_fragment)
31 print(nr)
```

8 Github + Useful links:

[GitHub Profile](#)

[Edit Distance - Youtube](#)

[The VJudge Problem - Edit Distance](#)

[Useful Python Guide - Youtube](#)

[Levenshtein Algorithm \(Explained In Details\)](#)

[StackOverflow](#)