Assignment 1

2016310932 배현웅

1. Overview

- A. Vigenere varint cipher Decryption is divided into two tasks. First task is to know the length of key. Second task is to distinguish each key value.
- B. To identify the length of key, I apply the statistics which alphabet distribution is not kind of uniform distribution. As the lecture note said, I use sigma of probabilities square. ($\sum p^2$)
- C. After we find the key length, we also check distribution of Nth letter. This time, I multiply frequency of alphabet included 'space' and the probability of each alphabet which I find distribution of Nth letter.

2. Code description

```
#include <stdio.h>
#include <string.h>
#include <stdlib.h>
#include <unistd.h>
#define SQUARE(x) ((x)*(x))
typedef struct twin{
    int max;
    int nu;
}twin;
char alpha[27] = {" etaonisrhlducmfwgpybvkxjqz"};
double percent[27] = {18.31, 12.60, 9.22, 8.01, 7.56, 7.00, 6.95, 6.52, 6.12,
    0.98, 0.69, 0.17, 0.12, 0.10, 0.06};
int length_of_key;
int main(){
    unsigned char ch;
    double max =0;
    length of key=0;
    double tmp_sum =0;
    FILE *fpIn, *fpOut;
    char buf[5001];
```

This code consists of initiation of value, such as length_of_key, percent and alpha which stands for the frequency of alphabet. Alpha[7] is corresponds to percent[7]. ('e' : 12.60% frequency)

Char buf value is for input text whose name is hw1_input.txt

```
fpIn = fopen("hw1_input.txt","r");
28
         fseek(fpIn,0,SEEK_END);
         int buf len = ftell(fpIn);
29
         fseek(fpIn,0,SEEK_SET);
         fread(buf,buf_len, sizeof(char), fpIn);
         fclose(fpIn);
         int count[260]={0, };
         memset(count,0,sizeof(int)*260);
         for(int i=1; i<=10; ++i){
             tmp_sum =0;
             memset(count,0,sizeof(int)*260);
             int j;
             for(j=1; i*j<buf_len; j++){</pre>
41
42
                  int tmp = (int)buf[i*j];
                  count[tmp]++;
44
             }
45
             for(int k=65; k<=90; ++k){</pre>
47
                  tmp_sum += SQUARE((double)count[k])/SQUARE((double)j);
             for(int k=97; k<=122; ++k){</pre>
49
                  tmp_sum += SQUARE((double)count[k])/SQUARE((double)j);
             if(max < tmp_sum){</pre>
                      max = tmp_sum;
                      length of key = i;
```

In line 27-32, I check file length by fseek and ftell incase of bumping into EOF value(in ASCII code 0x1A)

In line 36, to check possible length which is 1 to 10, I write for loop.

Line 41-44: traverse all buf element(input context), then save frequency of each asci code into count array.

Line 46-51 : square of probability of each alphabet (lower case : $65(a) \sim 90(z)$, capital : $97(A) \sim 122(Z)$ is added into tmp_sum

Line 52-54 : to compare value of each case (depends on key length) I take sigma of probabilities square. ($\sum p^2$) value and find the highest value among them. At that value, the highest possible key length is it.

```
61
         unsigned char key[length of key];
         memset(key,0,sizeof(unsigned char)*length_of key);
62
     FIND KEY:
64
         for(int k=0; k<length_of_key; ++k){</pre>
65
66
             int counts[260];
             memset(counts,0,260*sizeof(int));
67
68
             for(int i=0; length of key*i+k<=buf len; i++){</pre>
                 int pos = length_of key*i+k;
70
                 counts[buf[pos]] = counts[buf[pos]]+1;
71
72
73
             twin max_cnt,submax_cnt;
74
             max_cnt.max =0;
             submax cnt.max = 0;
75
```

After we find the key length, as I comment on code, we should find the key value

Line 63 is label for correction of key length. I explain it later.

Line 65-72: I take Nth letter in input text, and do similar thing as I find key length. However, this time I assume the repeated letter is 'space' which statistically occur often most

```
for(int i=0; i<=0xFF; ++i){</pre>
76
77
                  if(max_cnt.max < counts[i]){</pre>
78
                      submax_cnt.max = max_cnt.max;
79
                      submax_cnt.nu = max_cnt.nu;
80
                      max_cnt.max = counts[i];
81
                      max cnt.nu = i;
                  } else if(submax_cnt.max < counts[i]){</pre>
82
83
                      submax_cnt.max = counts[i];
                      submax cnt.nu = i;
84
85
86
              if(max_cnt.max*0.1 + submax_cnt.max > max_cnt.max){
87
88
                  if(length_of_key *2 <= 10){</pre>
89
                      length_of_key *=2;
90
91
                      goto FIND_KEY;
92
                  }
              int ke = max cnt.nu^' ';
94
95
              int ke1 = submax_cnt.nu^' ';
              key[k] = ke;
```

Line 76: we find most frequent letter(=max) and 2nd frequent letter(=submax) which are candidate

of 'space'

Line 87 – 93 : there is little difference b/w submax and max. There is unfolded key stream. So multiply 2, then do FIND_KEY work again (ex key : 0x01 0x02 0x01 0x03)

Line 94-95 save to ke and ke1 value which are candidate of key

```
97
              if(ke != ke1){
 98
                  double prob1 = 0;
99
                  double prob2 = 0;
                  prob1 = percent[0]*max_cnt.max;
100
                  prob2 = percent[0]*submax_cnt.max;
101
102
                  for(int i=1; i<27; ++i){
                      int diff = (int)(alpha[i]-' ');
103
104
                      int pos1 = (max_cnt.nu + diff)%128;
105
                      int pos2 = (submax_cnt.nu + diff)%128;
                      prob1 += percent[i] * counts[pos1];
106
107
                      prob2 += percent[i] * counts[pos2];
108
                  if(prob2 > prob1){
109
110
                      key[k] = ke1;
                  } //else printf("ke : %x\n",ke);
111
112
113
```

If two key candidate is different, I compare statistic values which is supposed to be ke or ke1is space

Compare sigma of alphabet probabilities square. ($\sum p^2$) which include 'space'

```
114
115
          for(int i=1; i<length_of_key; ++i){</pre>
116
               int same=0;
117
               if(i == 1){}
                   for(int j=0; j<length_of_key; j++){</pre>
118
119
                        if(key[j] == key[j+1]) same++;
120
                        else break;
121
122
                   if(same == (length of key-1)){
123
                        length_of_key = 1;
                        break; // exit the period test
124
125
                       .se same =0;
```

If there is repeated periodic value, we reduced to optimized way

```
(ex key : 0x01 0x02 0x01 0x02 -> 0x01 0x02 ))
```

So, I Check all possible value which is lower than length_of_key

Line 117-125: Test whether the key length is one

```
}else if(length of key%i == 0){
126
127
                   int period = length_of_key/i;
128
                   int okay =0;
                   if(key[0] == key[i]){
129
130
131
                       for(int iter=0; iter<(period-1); ++iter){</pre>
132
                           same=0;
133
                           for(int j=0; j<i; ++j){
134
                                if(key[iter*i+j] == key[(iter+1)*i+j]) same++;
                               else break;
135
136
                           if(same == i) okay++;
137
138
                           else break;
139
                       if((okay+1) == period) length_of_key = i;
140
141
                  }
142
143
              } // end else if(length_of_key%i == 0){
144
          } // end for(int i=1; i<length_of_key; ++i){
```

If key length is not one, then test possible key length which is aliquot(약수) of key length

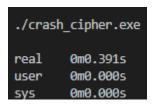
```
fpOut = fopen("hw1_output.txt","wb");
153
154
          // for file write
          fprintf(fpOut,"%#x",key[0]);
155
          for(int i=1; i<length_of_key; ++i){</pre>
156
               fprintf(fpOut, " %#.2x", key[i]);
157
158
          fprintf(fpOut, "\n");
159
160
          for(int i=0; i<buf_len; ++i){</pre>
161
              ch = buf[i]^key[i%length of key];
162
               fwrite(&ch, sizeof(ch), 1, fpOut);
163
164
          }
165
          fclose(fpOut);
166
          return 0;
168
```

After find the key value and key length, write down to output value.

3. Performance Check

I Test time command two text which is encrypted by key length 10 for performance check

- Text which consists of 1485 letters



- Text which consists of 4775 letters

```
real 0m0.043s
user 0m0.000s
sys 0m0.016s
```

Worst case for my code runtime is about 0.391 seconds