

Report of Assignment 7

Haodong Liao & Liyuan Zhang

July 12, 2019

1 Preface

This is a \LaTeX report written by Haodong Liao and Liyuan Zhang, students of UESTC who participated in the summer school of UCPH. More codes of this summer school are at this [repository](#).

2 Introduction

A string in F# is a sequence of characters¹, and some of the string functions are similar to list functions, after the practicing of the list function, it is time to move on.

In this assignment, we worked with simple text processing and met the requirements of analysing and generating some text according to statistics. In detail, we used pattern-matching and explicit type definition to create our own recursive functions, besides, we used powerful List-library functions to read, convert, analyze and generate target text.

3 Analysis and Design

3.1 Text Processing of character

3.1.1 ReadText

```
1 readText: filename:string -> string
2
3 // readFile.fsx
4 let filename = "readFile.fsx"
5 let line =
6     try
7         let reader = System.IO.File.OpenText filename
8         reader.ReadToEnd ()
9     with
10         _ -> "" // The file cannot be read, so we return an empty string
11
12 printfn "%s" line
```

In this sub-assignment, we needed to convert *readFile.fsx* into a more flexible function which reads the content of any text file.

Refer to the input processing code in assignment 6, there were two branches needed to be handled, one was to set the entrypoint of our code, the other was to change the fixed argument to a flexible one so that the function could read the content of any text file. The key part of the pseudocode is shown as following:

```
1 [<EntryPoint>]
2 let main (input argument) =
3     if (argument correct)
4         call readText function with responding name of file
5     else
6         display error message
```

3.1.2 ConvertText

```
1 convertText: src:string -> string
```

¹<https://docs.microsoft.com/en-us/dotnet/fsharp/language-reference/lists>

The goal of this sub-assignment was to convert letters of a string to lower case and removes all characters except a..z. According to [1], there was a *ToLower()* function *returns a copy of the string where each letter has been converted to lower case*, so we could use it to finish the 'convert' part and moved on to 'remove' part. The key part of the pseudocode is shown as following:

```
1 let convertText inputString =
2   let lowerStr = inputString.ToLower()
3   if (element of lowerStr was in the range of a..z)
4     keep the current element
5   else remove the current element
```

3.1.3 Get Histogram

```
1 histogram: src:string -> int list
```

We needed to counts occurrences of each lower-case letter of the English alphabet in a string and returns a list in this sub-assignment. According to [1], there was a *String.filter* function similar to *List.filter* function, we used it to filter the specific letter of the string and counts the length of it to get the occurrences. The key part of the pseudocode is shown as following:

```
1 let countchar (src:string)(letter:char) =
2   if letter was in the range of a..z then
3     filter the current letter
4     handle the rest string
5   else
6     returned an empty list
7 let histogram (src:string): int list = (countchar src 'a')
```

3.1.4 Generate Random String

```
1 randomString: hist: int list -> len:int -> string
```

The requirement of this sub-assignment was to change the given program and generate a string of a given length, and contains random characters distributed according to a given histogram which used our own functions. The flow of the program was shown as following:

```
1 step1 - Read the target text file and saved it
2 step2 - Converted the text to lower case & removed characters except a..z
3 step3 - Got the histogram of the text
4 step4 - Generated a new random string according to the histogram
```

3.2 Text Processing of pair of characters

3.2.1 Get Cooccurrence

In this sub-assignment, we needed to count occurrences of each pairs of lower-case letter of the English alphabet in a string and return a list of lists. Our solution was to use double recursion and process the letter in each pair separately. The key part of the pseudocode is shown as following:

```
1 // count the occurrence of each pair
2 let countnum =
3   if the length of src > 2 then
4     if the first two char = charpair then
5       1 + sum recursively
6     else
7       ignore current and sum recursively
8   else
9     return 0
10
11 // handle the second letter
12 let countcharpair =
```

```

13     if the second letter <= 'z' then
14         sum the occurrence and process the left string recursively
15     else return an empty list
16
17 // process the first letter
18 let countchar =
19     if the first letter <= 'z' then
20         process the second letter and process the left string recursively
21     else return an empty list
22
23 // call function from here
24 let cooccurrence (src:string) = countchar src 'a'

```

3.2.2 Markov Model

```

1 fstOrderMarkovModel: cooc: int list list -> len:int -> string

```

We needed to generate a random string of length *len*, whose character pairs were distributed according to a specified cooccurrence histogram *cooc* of the story "Little Claus and Big Claus" in this sub-assignment. The flow of the program was shown as following:

```

1 step1 - Count the occurrence of each pairs (like 'aa' to 'az')
2 step2 - Count the cooccurrence of all pairs ('aa' to 'zz')
3 step3 - Calculate the ratio of each pairs
4 step4 - Generate a random string of length len
5         according to a given cooccurrence and the ratio of each pairs

```

4 Program description

Our implementation of **readText** function was as follows:

```

1 // 7.1
2 let printErrorMessage () =
3     printfn "Program Input should be 'readfile filename'"
4
5 let readText (filename:string) : string =
6     let line =
7         try
8             let reader = System.IO.File.OpenText filename
9             reader.ReadToEnd ()
10         with
11             _ -> "" // The file cannot be read, so we return an empty string
12     printfn "%s" line
13     line
14
15 [<EntryPoint>]
16 let main (paramList : string[]) : int =
17     if paramList.Length <> 2 then
18         printErrorMessage ()
19         0
20     else
21         match paramList.[0] with
22         | "readfile" ->
23             let str = readText paramList.[1]
24             printfn "str : %A" str
25         | _ ->
26             printErrorMessage ()
27     1

```

As it shows, the key of this function was to match the argument list, the first argument should be the command 'readfile', and the second argument was the name of file.

Implementation of **convertText** function was as follows:

```

1 // 7.2
2 let rec convertText (src:string) : string =
3     let lowerSRC = src.ToLower()
4     match lowerSRC with
5     | "" -> ""
6     | elm ->
7         if ('a' <= elm.[0] && elm.[0] <= 'z')
8         then
9             elm.[0].ToString() + (convertText elm.[1..])
10        else
11            convertText elm.[1..]

```

It made us wondering whether there was an expression of string could work as 'elm::rest' to represent a list at the beginning, things did not go well so we changed our mind to use elm[1..] to represent the rest of a string.

Implementation of **histogram** function was as follows:

```

1 // 7.3
2 let rec countchar (src:string)(letter:char): int list =
3     if letter >= 'a' && letter <= 'z' then
4         ((String.filter (fun char -> char=letter) src).Length)::
5         (countchar src (char (int letter + 1)))
6     else []
7 let histogram (src:string): int list = (countchar src 'a')

```

This function was the key to solving the following problems.

Assignment 7.4 was a summary of the previous three assignments, and the key part of implementation of it was as follows:

```

1 /// step1 readfile
2 let str = readText "littleClausAndBigClaus.txt"
3 /// step2 convertfile
4 let convertRes = convertText str
5 /// step3 get histogram
6 let hist = histogram convertRes
7 let alphabet = List.init hist.Length (fun i -> 'a' + char i)
8 printfn "A histogram:\n %A" (List.zip alphabet hist)
9 /// 7.4 step4 generages a new random string
10 let ranStr = randomString hist convertRes.Length
11 printfn "A random string: %s" ranStr
12 let newHist = histogram ranStr
13 printfn "Resulting histogram:\n %A" (List.zip alphabet newHist)

```

It was important to have a clear thought of processing flow.

Implementation of **cooccurrence** function was as follows:

```

1 // 7.5
2 let rec countnum (src:string) (charpair:string):int =
3     if src.Length >=2 then
4         if src.[0..1] = charpair then
5             1+(countnum src.[1..] charpair)
6         else
7             (countnum src.[1..] charpair)
8     else
9         0
10 let rec countchpair (src:string) (letter1:char) (letter2:char):int list =
11     if letter2 <= 'z' then
12         let charpair = string letter1 + string letter2
13         (countnum src charpair)::(countchpair src letter1 (letter2+char 1))
14     else []
15
16 let rec countchar (src:string) (letter1:char):int list list =
17     if letter1 <= 'z' then
18         (countchpair src letter1 'a')::(countchar src (letter1 + char 1))
19     else []

```

```

20 let cooccurrence (src:string):int list list =
21     countchar src 'a'

```

This was a difficult problem for us, and the key to the solution was the idea of 'divide and conquer', if we could not process the pair of letters at a time, we handle the letter of it separately.

Assignment 7.6 was a summary of all the previous assignments, and the key part of implementation of it was as follows:

```

1  let fstOrderMarkovModel (cooc:int list list) (len:int):string =
2      let strLengthlist =
3          List.map (fun x -> (float x)) (List.map (List.sum) cooc)
4      let strLengthSum = List.sum strLengthlist
5      let ratiolist = List.map (fun x -> x/strLengthSum) strLengthlist
6      let randomStringAlt (hist: int list):string = randomString hist
7          (int((float (List.sum (hist))/strLengthSum)*(float len)))
8      let strlist = List.map (randomStringAlt) cooc
9      let composestr (acc:string) (elm:string):string = acc + elm
10     List.fold composestr "" strlist
11     printfn "test: \n %A" (fstOrderMarkovModel (cooccurrence a) (a.Length))
12     printfn "present cooc : \n %A"
13     (cooccurrence (fstOrderMarkovModel (cooccurrence a) (a.Length)))

```

This was another difficult assignment for us, the main challenge was to arrange the text properly so the occurrence was similar. We tried to generate a new character according to the letter before the current letter, but things just did not worked well for reason we did not figured out, then we tried with the weight of the pair based on the proportion of the sum of the single letters to the sum of the occurrences of all the letters, for example: $Weight_{a'} = \frac{Occurrenceof'aa'to'az'}{Occurrenceof'aa'to'zz'}$. In this way, when there was a given length, the length of each character list (like 'aa'..'az') would be as the weight of all character list, so that we could generate a random string which had a similar distribution but different content.

5 Evaluation

The testing environment was macOS Mojave 10.14.5 system with iTerm and Microsoft (R) F# Compiler version 4.1.

We combined all the assignment into a single program, compiled and tested it with parameter *readfile littleClausAnd-BigClaus.txt*, and the key part of result (because it was to long) were showed as following:

```

→ fsharpc Assig2.fsx && mono Assig2.exe readfile littleClausAndBigClaus.txt
Microsoft (R) F# Compiler version 4.1
Copyright (c) Microsoft Corporation. All Rights Reserved.
str :
In a village there lived two men who had the self-same name. Both were named Claus. But one of t
hem owned four horses, and the other owned only one horse; so to distinguish between them people
called the man who had four horses Big Claus, and the man who had only one horse Little Claus.
Now I'll tell you what happened to these two, for this is a true story.

```

Figure 1: Testing of readText function

```

convertedstr :
"ina village therelivedtwomenwhohadtheselfsamenamebothwerenamedclausbutoneofthemownedfourhorsesand
theotherownedonlyonehorsesotodistinguishbetweenthepeoplecalledthemanwhohadfourhorsesbigclausand
themanwhohadonlyonehorselittleclausnowilltellyouwhathappenedtothesetwoforthisisatruestorythewhol
eweekthroughlittleclausadtoplowforbigclausandlendhimhisonlyhorseinreturnbigclauslenthimallfour
fhishorsesbutonlyforonedayaweekandthathadtobesundayeachsundayhowproudlylittleclauscrackedhiswhip
overallthefivehorseswhichwereasgoodashisownonthatdayhowbrightlythesunshonehowmerrywerethechurchb
ellsthatranginthesteepelhowwelldressedwereallthepeoplewhopassedhimwithhymnbookstuckedundertheira
rmsandastheywenttheirwaytochurchtoheartheparsonpreachhowthepeopleliddidstaretoeelittleclausplowing
withallfivehorsesthismadehimfeelsoproudthathewouldcrackhiswhipandhollogetupallmyhorsesyoumustnot

```

Figure 2: Testing of convertText function

6 Conclusion

In this assignment, I implemented the *myFold* and *myFilter* functions with my own recursive method. I stuck at the beginning of my writing, but things went smoothly when I rearranged my thoughts and wrote down the pseudocode. It's


```

histogram :
[1293; 253; 391; 810; 2036; 326; 319; 1308; 999; 28; 179; 821; 391; 865; 1262;
 207; 9; 785; 1000; 1591; 569; 125; 409; 24; 311; 5]
new random string1 :
"vltensrhtnueetdfsnhraattntjblgcatholuetmseditannotuaiftidworfotmiuhteohmrtutuhaletiohanmiehwmc
edhnwurrialhfnhalrudkieedhhtiwehyonsigictlihcaryehooahotttdinnteksdhoannefthvatcdohoeohnysthtllu
euowoaynsigeaosaagteeeedssoeohihhpaedknhaesfdmfvtfmveparlethrumdttoyrtttfueyhtoehftwslcpohruoen
croietgeahpoflagomnaibnghswfhwcwecgerhhabtmctseoaaahadadhcluatuilmmyhrngeewheaekysohanetebdhetia

```

Figure 3: Testing of histogram and randomString function

```

new histogram :
[1328; 250; 377; 806; 2102; 341; 310; 1321; 983; 29; 195; 825; 393; 866; 1244;
 189; 9; 782; 987; 1588; 516; 136; 401; 27; 310; 1]
str's cooc:
[[0; 35; 64; 77; 0; 12; 31; 7; 84; 0; 26; 64; 28; 228; 0; 16; 0; 100; 147; 141;
 122; 36; 26; 3; 43; 3];
 [19; 0; 0; 0; 64; 0; 0; 0; 44; 0; 0; 12; 0; 0; 29; 0; 0; 8; 0; 0; 61; 0; 0; 0;
 16; 0];
 [70; 0; 0; 0; 21; 0; 0; 54; 2; 0; 60; 112; 0; 0; 50; 0; 0; 15; 0; 3; 4; 0; 0; 0;
 0; 0];

```

Figure 4: Histogram and Cooc of Generate String of randomString function

```

new random string2 :
"tndrplctngcdudunryistnnlutnntvnuiltfnwnsttblslnnsttncsmnuttsurbtrtiusmtbrsruvnisrnwtwgutswni
ristitnnntdvtitnsnctcimvtcnunlswnnnztsgusmguugdlurlppymuntcswwdggnutticigtctcuvntgndsyuns
uinucncnsntdlrumgntvrtrtstttznndttmynmnglpsrtpndndssikrmniblsgtcpunnuknhsnlaiiywnbbsdh
dtutuinusnunlsztrusnyuwnfkicvsgcwnstrhnkttsvntvlwrndbsiltnrswmwyitpvyiknnnsifbcnmvdnibbmgnbd
stlntttriyltncsrwvgnflcistititrigttginnlrinsnbinnilltnnikivrtbitrikvnsrtnuvwnusrtmtrmvnniv

```

Figure 5: Testing of fstOrderMarkovModel function

```

new random string2's cooc :
[[184; 12; 35; 47; 218; 32; 15; 108; 91; 2; 5; 68; 28; 32; 117; 10; 1; 44; 63;
 96; 29; 5; 31; 1; 29; 0];
 [21; 6; 6; 11; 33; 2; 2; 23; 21; 1; 2; 12; 8; 10; 16; 2; 0; 15; 14; 16; 10; 2;
 6; 0; 5; 0];
 [27; 9; 26; 34; 20; 11; 11; 27; 14; 0; 2; 16; 10; 29; 22; 3; 1; 35; 32; 39; 8;
 4; 13; 1; 6; 0];
 [49; 17; 23; 96; 52; 25; 54; 35; 38; 3; 6; 53; 15; 59; 37; 10; 1; 51; 66; 116;
 17; 10; 17; 4; 9; 0];

```

Figure 6: Cooccurrence of fstOrderMarkovModel function

easy to get bogged down in the details of a program, but we should take a top-down functional programming approach and thinking more about what to do than how to do it.

References

[1] Jon Spurring. *Learning to Program with F#*. 2019.