

Text Generation using a Markov Chain ^{★, ★★}

Topics

- Strings, Lists, Hash tables, Arrays
- A little bit of Format and Arg

Introduction

The goal of this exercise suite is to decompose an existing text corpus, in order to recompose sentences by randomly collating original word sequences with the same appearance probability.

For this, we will use a markov chain style method. We will construct a markov chain from an input text, whose states are the words, and transitions are quantified by the probability for a word to be the predecessor of another.

For instance, if we examine *"I am a man and my dog is a good dog and a good dog makes a good man."*, delimiting it with "START" and "STOP" to identify proper beginning and end of sentences, we end up with the following chain, that can then be walked to generate new text that resembles the input.

"START" → "I" (100%)	"my" → "dog" (100%)
"I" → "am" (100%)	"dog" → "is" (33%), "and" (33%), "makes" (34%)
"am" → "a" (100%)	"good" → "dog" (66%), "man" (34%)
"a" → "man" (25%), "good" (75%)	"is" → "a" (100%)
"man" → "and" (50%), "STOP" (50%)	"makes" → "a" (100%)
"and" → "my" (33%), "a" (67%)	

Exercise 1 – Basic Version

Our first goal will be to build such a table and generate sentences from it, quick and dirty style, using lists and their predefined operators. Below is an example output.

```
1 0: START I am a good man STOP
2 1: START I am a good dog is a good dog and my dog and my dog is a man and
3   my dog and a man STOP
4 2: START I am a good dog is a man and my dog makes a good man STOP
5 3: START I am a good dog makes a good dog is a good dog and a good dog
6   makes a good dog is a man STOP
7 4: START I am a good dog and a man and a good dog and a good man and a good
8   dog is a good dog is a good man and a man STOP
9 5: START I am a good dog and a good dog and my dog is a man STOP
10 6: START I am a man STOP
11 7: START I am a good dog is a good dog is a good dog and my dog is a man
12   STOP
13 8: START I am a good man STOP
14 9: START I am a good dog makes a good dog and a good dog is a good dog is a
```

```

15     good man and my dog is a good dog and my dog and a good man and a good
16     dog is a good man STOP
17 10: START I am a man and my dog and my dog is a good dog and a good dog
18     makes a man and my dog is a good man and my dog and my dog and my dog
19     makes a man STOP

```

Question 1.1 – Write a function `words_of_file : string -> string list` that takes a file name and returns the list of its words. As a first approximation, you can consider roman letter sequences as words and everything else as separators.

Question 1.2 – Write `build_table : string list -> (string * string list) list` that associates each word present in the input text to all its possible successors (so that we can access the successors using `List.assoc`). If a successor appears several times, simply add all the duplicates, so that we respect the probability distribution implicitly.

Question 1.3 – Write the random generation function `walk : (string * string list) list -> string -> string -> string list` which takes a table, a given starting word and stops the generation when a given final word is encountered. We assume that the request is well formed for the table, meaning that there is always a path in the table between the given start and stop.

Question 1.4 – Write a small main program that takes an input text file as parameter and generates and prints 20 sentences. In order to ensure consistency, wrap the input text with fake tokens "START" and "STOP", that you also pass to the walk function.

You can use the almighty `Format` module to pretty print the text as in the example output.

Exercise 2 – Performance Improvements

The previous prototype was quick to write, but does not scale. Let us improve it so that we can index a full book, or even several ones.

Question 2.1 – If we want to input a corpus bigger than a few dozens of lines, we have to build the table with a more appropriate data structure than an associative list. Update the `build_table` function so it uses a hash table internally, without changing its type. You can use `Hashtbl.fold` to dump the internal table to a list. Experience the performance boost by running the generator before and after on some big ebook.

Question 2.2 – Now, generation should be the bottleneck. You can enhance the performance doing two things: (1) dropping the associative list completely and (2) using arrays for successors instead of lists.

In order to keep all those things hidden, and to leave space for potential future improvements, we will hide all the implementation behind a function type.

Update `build_table` so that it returns a lookup function of type `(string -> string)` that does the searching and the random choosing.

Question 2.3 – Update `walk` accordingly.

Exercise 3 – Quality Improvements

When given books as input, our generator generates a lot of garbage. We will now enhance it to use the structure of its input in order to filter the results.

Question 3.1 – First, let us improve our trivial input splitting function. Update `words_of_file` to return a list of sentences instead of a single big succession of words.

A good first approximation is to consider end-of-word dots and skipped lines as sentence terminators, and capitals just following dots as sentence starters.

Question 3.2 – In order to modify as few things as possible, simply collate all sentences using the "START" and "STOP" tokens to separate them, before calling `build_table`.

Question 3.3 – Another way to enhance the output is to help our brain filter the results. An idea is to colorize the words with respect to the input text they have been extracted from.

Update the main program to take a list of inputs, and instead of simply building a big list of words to pass to `build_table`, build a list of pairs, associating each word to a number corresponding to the order of its originating source on the command line. You can associate files starting from 1, reserving 0 for the fake "START" and "STOP" tokens.

Question 3.4 – Update `build_table` so that the returned lookup function returns a pair of the selected successor word and its origination. The selection in the table must not take in consideration the color, so the lookup function still only takes a string as input.

Question 3.5 – Use these numbers in the main program to print ANSI escape sequences. In OCaml, escaped character codes are in decimal, so you must use `\027` to start escape sequences.

Exercise 4 – Prefix Matching

A drastical way of improving the results is, instead of considering the possible successors of each word, to consider the ones of each prefix of a given length N .

These examples have been generated from a mix of Arsène Lupin, a collection of sandwich recipes, and the OCaml Manual, using this algorithm with a prefix length set to 2.

- *And then, was she impressed with the g option to ocamlc during the past few years.*
- *No sound could be freely redistributed.*
- *The file currently being parsed by the compiler will be your humble servant.*
- *Devanne exclaimed with much gusto: Ah monsieur, you can observe the entire loaf, butter each slice a layer of the greatest integer less than or equal than the usual slices for sandwiches.*
- *The buttered side of the constrained signature.*
- *Dead Dead he repeated, to himself; they are carefully cooked, rejecting the skins; work this together.*
- *Sherlock Holmes is the source code for a moment.*
- *Readers unfamiliar with lex and yacc can report that expression of intense agony.*

Question 4.1 – Write `split_prefix : int -> (string * int) list -> string list * (string * int)` That takes an input list of colored tokens, and returns the uncolored prefix of the given length N and its successor token. The function should raise `Not_found` when the end of list is reached.

Question 4.2 – Update `build_table` to take the prefix length N and return a lookup function that takes a prefix of type `string list` as input. Use `split_prefix` to deconstruct the list until it raises `Not_found`.

Question 4.3 – Update `walk` so that it takes an input prefix sequence of length N instead of a single word, and stops when the accumulated sentence ends with a stop sequence of length N .

Question 4.4 – Instead of *"START"* and *"STOP"*, we will simply use repeated sequences of length N of these same tokens. Write a utility `repeat n elt` that builds a list containing n times `elt`, and update the main program.

Question 4.5 – Use the `Arg` module to parse a `-length` option to set the prefix length.

Text Generation using a Markov Chain

SOLUTIONS – SOLUTIONS – SOLUTIONS – SOLUTIONS

Solution to question 1.1

One in functional style:

```
1 let words_of_file fname =
2   let isletter c = match c with
3     | '0'-'9' | 'a'..'z' | 'A'..'Z' | '\128'..'255' -> true
4     | _ -> false in
5   let rec read fp acc =
6     match input_line fp with
7     | line ->
8       let acc = split line acc 0 0 in
9       read fp acc
10    | exception End_of_file ->
11      List.rev acc
12  and split line acc s i =
13    if i < String.length line && isletter line.[i] then
14      split line acc s (succ i)
15    else if s < String.length line then
16      let acc =
17        if s <> i then
18          String.sub line s (i - s) :: acc
19        else acc in
20      split line acc (succ i) (succ i)
21    else acc in
22  let fp = open_in fname in
23  let words = read fp [] in
24  close_in fp ;
25  words
```

One in imperative style:

```
1 let words_of_file fname =
2   let fp = open_in fname in
3   let words = ref [] in
4   (try
5     while true do
6       let line = input_line fp in
7       let rec eat s i =
8         if i < String.length line then
9           match line.[i] with
10            | '0'..'9' | 'a'..'z' | 'A'..'Z' | '\128'..'255' -> eat s (succ i)
11            | _ -> cut s i
12          else cut s i
13        and cut s i =
14          if s < String.length line then begin
15            if s <> i then
```

```

16         words := String.sub line s (i - s) :: !words ;
17         eat (succ i) (succ i)
18     end
19     in eat 0 0
20 done
21 with End_of_file -> () ;
22 close_in fp ;
23 List.rev !words

```

Solution to question 1.2

```

1 let build_table phrase =
2   let update table word next =
3     match table with
4     | [] -> [ (word, [ next ]) ]
5     | (w, l) :: ws ->
6       if w = word then
7         (w, next :: l) :: ws
8       else
9         (w, l) :: update ws word next in
10  let rec build phrase =
11    match phrase with
12    | [] | [ _ ] -> []
13    | word :: (next :: _ as rest) ->
14      let table = build rest in
15      update table word next in
16  build phrase

```

Solution to question 1.3

```

1 let walk table start stop =
2   let rec walk index =
3     if index = stop then
4       [ stop ]
5     else
6       let choices = List.assoc index table in
7       let next = List.nth choices (Random.int (List.length choices)) in
8       index :: walk next in
9   walk start

```

Solution to question 1.4

```

1 let () =
2   Random.self_init () ;
3   match Array.to_list Sys.argv with
4   | [ exe ; file ] ->
5     let start = "START" and stop = "STOP" in
6     let words = [ start ] @ words_of_file file @ [ stop ] in
7     let table = build_table words in
8     for i = 0 to 20 do
9       let sentence = walk table start stop in
10      Format.printf "%_3d:_@" i ;
11      List.iter (Format.printf "%s@_") sentence ;
12      Format.printf "@]@."
13   done

```

```

14 | exe :: _ ->
15   Format.printf "Usage:_%s_<input_text_file>\n" exe ;
16   exit 1
17 | [] ->
18   exit 1

```

Solution to question 2.1

```

1 let build_table phrase =
2   let table = Hashtbl.create 1000 in
3   let update word next =
4     match Hashtbl.find table word with
5     | succs -> succs := next :: !succs
6     | exception Not_found -> Hashtbl.add table word (ref [ next ]) in
7   let rec build phrase =
8     match phrase with
9     | [] | [ _ ] -> ()
10    | word :: (next :: _ as rest) ->
11      update word next ;
12      build rest in
13   build phrase ;
14   Hashtbl.fold (fun word succs r -> (word, !succs) :: r) table []

```

Solution to question 2.2

```

1 let build_table phrase =
2   let table = Hashtbl.create 1000 in
3   let update word next =
4     match Hashtbl.find table word with
5     | succs -> succs := next :: !succs
6     | exception Not_found -> Hashtbl.add table word (ref [ next ]) in
7   let rec build phrase =
8     match phrase with
9     | [] | [ _ ] -> ()
10    | word :: (next :: _ as rest) ->
11      update word next ;
12      build rest in
13   build phrase ;
14   let final_table = Hashtbl.create 1000 in
15   Hashtbl.iter
16     (fun word succs -> Hashtbl.add final_table word (Array.of_list !succs))
17     table ;
18   let lookup word =
19     let succs = Hashtbl.find final_table word in
20     succs.(Random.int (Array.length succs)) in
21   lookup

```

Solution to question 2.3

```

1 let walk table start stop =
2   let rec walk index =
3     if index = stop then
4       [ stop ]
5     else
6       let next = table index in

```

```

7     index :: walk next in
8     walk start

```

Solution to question 3.1

```

1  let words_of_file fname =
2      let isletter c = match c with
3          | '0'..'9' | 'a'..'z' | 'A'..'Z' | '\128'..'255'
4          | '.' | '\' | ';' | ',' | '-' | '_' | ':'
5          | '(' | ')' | '[' | ']' | '{' | '}' -> true
6          | _ -> false in
7      let empty =
8          ([], []) in
9      let grab (cur, acc) =
10         if cur == [] then acc else List.rev cur :: acc in
11      let push word (cur, acc) =
12         match word.[0], word.[String.length word - 1], cur with
13         | _, '.', _ -> (cur, grab (word :: cur, acc))
14         | 'A'..'Z', _, [] -> ([ word ], acc)
15         | _, _, [] -> (cur, acc)
16         | _ -> (word :: cur, acc) in
17      let rec read fp acc =
18         match input_line fp with
19         | line ->
20             let acc = split line acc 0 0 in
21             read fp acc
22         | exception End_of_file ->
23             grab acc
24      and split line acc s i =
25         if i < String.length line && isletter line.[i] then
26             split line acc s (succ i)
27         else if s < String.length line then
28             let acc =
29                 if s <> i then
30                     push (String.sub line s (i - s)) acc
31                 else acc in
32             split line acc (succ i) (succ i)
33         else acc in
34      let fp = open_in fname in
35      let sentences = read fp empty in
36      close_in fp ;
37      sentences

```

Solution to question 3.2

```

1  ...
2  | [ exe ; file ] ->
3      let start = "START" and stop = "STOP" in
4      let words =
5          List.fold_right
6              (fun words acc -> [ start ] @ words @ [ stop ] @ acc)
7              (words_of_file file) [] in
8      let table = build_table words in
9      ...

```


Solution to question 3.3

```
1 ...
2 | exe :: files ->
3   let start = "START", 0 and stop = "STOP", 0 in
4   let words =
5     List.mapi (fun i file ->
6       List.fold_right
7         (fun words acc ->
8           let words = List.map (fun w -> w, succ i) words in
9           [ start ] @ words @ [ stop ] @ acc)
10      (words_of_file file) [])
11   files
12 |> List.flatten in
13 let table = build_table words in
14 ...
```

Solution to question 3.4

```
1 let walk table start stop =
2   let rec walk index =
3     if index = stop then
4       [ stop ]
5     else
6       let next = table (fst index) in
7       index :: walk next in
8   walk start
```

Solution to question 3.5

```
1 ...
2   for i = 0 to 20 do
3     let sentence = walk table start stop in
4     Format.printf "%_3d:_@" i ;
5     List.iter (fun (word, color) ->
6       let color_es = Format.asprintf "\027[%dm" (30 + color) in
7       let reset_es = "\027[0m" in
8       Format.printf "@<0>%s%s@<0>%s@_" color_es word reset_es)
9     sentence ;
10    Format.printf "@]@."
11  done
12 ...
```

Solution to question 4.1

```
1 let split_prefix length l =
2   let rec split_prefix n acc l = match (n, l) with
3   | 0, next :: _ -> List.rev acc, next
4   | _, (next, _) :: rest -> split_prefix (pred n) (next :: acc) rest
5   | _ -> raise Not_found
6   in split_prefix length [] l
```

Solution to question 4.2

```
1 let build_table length phrase =
2   let table = Hashtbl.create 1000 in
```

```

3  let update word next =
4      match Hashtbl.find table word with
5      | succs -> succs := next :: !succs
6      | exception Not_found -> Hashtbl.add table word (ref [ next ]) in
7  let rec build phrase =
8      match phrase with
9      | [] | [ _ ] -> ()
10     | _ :: rest ->
11         match split_prefix length phrase with
12         | prefix, next ->
13             update prefix next ;
14             build rest
15         | exception Not_found -> () in
16  build phrase ;
17  let final_table = Hashtbl.create 1000 in
18  Hashtbl.iter
19      (fun word succs -> Hashtbl.add final_table word (Array.of_list !succs))
20      table ;
21  let lookup word =
22      let succs = Hashtbl.find final_table word in
23      succs.(Random.int (Array.length succs)) in
24  lookup

```

Solution to question 4.3

```

1  let walk table start stop =
2      let roll l w =
3          List.tl l @ [ w ] in
4      let rec walk index =
5          if index = stop then
6              []
7          else
8              let next = table (List.map fst index) in
9              next :: walk (roll index next) in
10     start @ walk start

```

Solution to question 4.4

```

1  let () =
2      Random.self_init () ;
3      let length = 2 in
4      match Array.to_list Sys.argv with
5      | [ exe ] ->
6          Format.printf "Usage:_%s_<input_text_file>\n" exe ;
7          exit 1
8      | exe :: files ->
9          let rec repeat n elt =
10              if n = 0 then [] else elt :: repeat (pred n) elt in
11          let start = repeat length ("START", 0)
12          and stop = repeat length ("STOP", 0) in
13          let words =
14              List.mapi (fun i file ->
15                  List.fold_right
16                      (fun words acc ->
17                          let words = List.map (fun w -> w, succ i) words in

```

```

18         start @ words @ stop @ acc)
19         (words_of_file file) [])
20     files
21     |> List.flatten in
22     let table = build_table length words in
23     ...

```

Solution to question 4.5

```

1  let () =
2      Random.self_init () ;
3      let length = ref 2 in
4      let rem = ref [] in
5      let exe = Filename.basename Sys.executable_name in
6      let usage = exe ^ "[options]_<text_file>_<text_file>_..." in
7      Arg.parse
8          [ "-length", Arg.Set_int length, "set_the_length_of_the_prefix" ]
9          (fun arg -> rem := arg :: !rem)
10         usage;
11     match List.rev !rem with
12     | [] ->
13         Format.printf "Usage:_%s\n" usage ;
14         exit 1
15     | files ->
16         let rec repeat n elt =
17             if n = 0 then [] else elt :: repeat (pred n) elt in
18         let start = repeat !length ("START", 0)
19         and stop = repeat !length ("STOP", 0 )in
20         let words =
21             List.mapi (fun i file ->
22                 List.fold_right
23                     (fun words acc ->
24                         let words = List.map (fun w -> w, succ i) words in
25                         start @ words @ stop @ acc)
26                     (words_of_file file) [])
27             files
28         |> List.flatten in
29         let table = build_table !length words in
30         for i = 0 to 20 do
31             let sentence = walk table start stop in
32             Format.printf "%_3d:_@[" i ;
33             List.iter (fun (word, color) ->
34                 if color <> 0 then
35                     let color_es = Format.asprintf "\027[%dm" (30 + color) in
36                     let reset_es = "\027[0m" in
37                     Format.printf "@<0>%s%s@<0>%s@_" color_es word reset_es)
38                 sentence ;
39             Format.printf "@]@."
40         done

```