LKR Project 3

LKR Page 2 on 17

Table des matières

1	Structural Information Theory Grammar	3
2	Most Frequent Consecutive Repeated Sub-string	5
3	Symmetry matching	10
4	Generating the final factorization	11
5	Other approaches we tried	13
	5.1 With a DCG Grammar	13
	5.2 Alternative to our current iteration algorithm	14
6	Program Usage	15
7	Shortlog of the teamwork	16

LKR Page 3 on 17

1 Structural Information Theory Grammar

By studying the Structural Information Theory, we came up with the following DCG grammar for verifying it:

```
1 exp --> term.
 2 \exp --> term, [+], exp.
3
4 term --> factor.
5
  term --> digit, [*], exp.
6
7 factor --> elem.
  factor --> ['S'], ['['], sym, [']']. %S[(A)(B),(C)]
  factor --> ['<'], alt, ['>'], ['/'], ['<'], alt, ['>'
      ]. %<(A)>/<(B)(C)(D)>
10 factor --> ['('], exp, [')'].
11
12 \text{ sym } --> \text{ factor.}
13 sym --> factor, [','], factor.
14 sym --> factor, sym.
15
16 alt --> factor.
17
   alt --> factor, alt.
18
19 elem --> char.
20 elem --> char, elem.
21 char --> [D], {is_alnum(D)}.
22 digit --> [D], {is_alnum(D)}.
23 digit --> [D], {number(D)}.
```

LKR Page 4 on 17

The grammar in itself is not very complex, but we didn't use it to factorize, it was just used as an output verification step.

Instead, we implemented several algorithms to match the patterns and then, we added a step to consume characters and generate the factorization. This will be explained later on.

LKR Page 5 on 17

2 Most Frequent Consecutive Repeated Substring

The iteration rule from standard information theory aims to factorize the consecutive repeated pattern of a given string. Here we have some examples:

- -- AAAAA <= 6*(A)
- -- AXDEXDEXDEX <= AX + 3*(DEX)
- -- DEXDEXDEXAX <= 3*(DEX) + AX
- OKDEXDEXAX \leq OK + 3*(DEX) + AX

To implement this feature, we applied the Longest Repeated Substring Algorithm found in the following link plus an index checking to search only for **consecutive** patterns.

LKR Page 6 on 17

Theses are the steps of this algorithm:

1. First we have to build a suffix array.

We 'll use "abracadabra for this example".

index: word

0: abraabraabra

1: braabraabra

2 : raabraabra

3: aabraabra

4: abraabra

5: braabra

6 : raabra

7: aabra

8 : abra

9 : bra

10 : ra

11:a

2. Then we iteration through the array to compare all suffixes to each other. We take a suffixes and we check if it's a prefixing the others.

```
prefix(abra, abraabra) is true
prefix(abra, abraabra) is true
prefix(abra, abra) is true
The rest is false.
```

LKR Page 7 on 17

3. But here we want to find only consecutives patterns:

```
"abraabra" is containing a consecutive pattern "abra".
```

"abraxabra" is a not containing any consecutives patterns.

"abrabra" contains two "abra" pattern but they're overlaping so we want to get rid of it.

We can validate the consecutivity of a pattern thanks to a index checking:

```
1
      is_continuous(_, SubIdx, SuffIdx, _, _) :-
2
      SubIdx == SuffIdx,
3
4
      is_continuous(SubLen, SubIdx, SuffIdx, Sub,
          Suff) :-
5
      drop(SubLen, Suff, NewSuff),
6
      prefix(Sub, NewSuff),
7
      NewSuffIdx is SuffIdx + SubLen,
      is_continuous(SubLen, SubIdx, NewSuffIdx,
8
          Sub, NewSuff).
```

LKR Page 8 on 17

Comparing "abra" of index 8 and "abraabra", of index 4, first check if abra is prefixing abraabra, prefix(abra, abraabra) = true, then drop as many letters than the length of "abra" in the suffix, drop(4, "abraabra", R). R = abra and increment his index, 4 + 4 = 8. The pattern is consecutive if the incermented index reach "abra"s index. Which is 8 so its valid.

Let's see for an overlaping example "abrabra":

"abra" is prefixing "abrabra", "abrabra"'s index is 0 and "abra"'s index is 3. When we drop(4, abrabra, R) R="bra", and the incremented index is 4 which is not equal to 3. We should continue untill prefix(abra, bra) is true but this case stops here.

- 4. To find the most frequent consecutive pattern, we just count the number of times the routine explained above succeed.
- 5. But there's still some cases failing, like "cozcozcozw".

 The reason is that we don't find "coz" in the suffix array. The solution is that for each suffix, we have to compare their own prefixes. More precisely, when we iteration on "cozw", we first compare it to all suffix array's elements, then drop the last char and re-compare again:

prefix(cozw, cozcozw) is false,

drop_last_char(cozw, R). R=coz.

We compute again with "coz" and we find 3*(coz).

LKR Page 9 on 17

6. At this point, we can find the most frequent consecutive repea-

ted pattern:

iteration("abraabra", Cnt, R). Cnt = 2, R = "abra".

But we want this result:

iteration("xabraabray", R). R= "x+2*(abra)+y".

We can reach this result with a simple split string separating the repeated substring and the rest then concatenate between "+" chars.

LKR Page 10 on 17

3 Symmetry matching

To find any symmetry pattern, we extended the suffix array generation algorithm to generate a prefix array as well.

Using this prefix array, we're going to be able to match the suffixes and prefixes the following way to find our patterns:

- 1. Generate all possible **consecutive** sub-lists of the input array for which we want to find the pattern.
- 2. Match a prefix and suffix that is equal to one of the possible **consecutive** sub-lists.

If we find a match, this means that we have a symmetry pattern within our given input string.

So, what we do, is we use the findall/3 function, we can get the whole set of possible symmetry patterns in our given input string and then, we choose the longest one.

After this, we inject the result in a function that will cut it into pieces, to be usable as an output.

For example:

```
adacbbcada yields S[(bb), (adac)].

cococo yields c+S[(o), (oc)].
```

LKR Page 11 on 17

4 Generating the final factorization

To be able to output a result, we probe the input list and in a recursive fashion, we check for each patterns and output them in a result list in the correct order.

Of course, we also need an heuristic to choose which factorization pattern to output over another.

Let's say a symmetry and an iteration can be seen in a string at the same place, for example :

```
oooooo can be seen as 7 * (0) or S[(0), (000)], or S[(0), (3 * (0))], or even S[(0), (S[(0), (00)])].
```

For this cases, we take each outputs generated by all conflicting patterns, and we compare their lengths stripped of special characters.

For instance, for an iteration, the stripped characters would be : $n\star ()$.

For the previous example, this would give :

```
Len ("7*(0)") < Len ("S[(0), (000)]") or 1 < 4 which means that the iteration patterns will be chosen over the symmetry.
```

After, we factorize this in a recursive fashion and stop when there is no change anymore.

Finally, this will yield a good factorization of our input string.

LKR Page 12 on 17

This also means that we can factorize an expression that is only semi-factorized.

LKR Page 13 on 17

5 Other approaches we tried

5.1 With a DCG Grammar

As explained above, we didn't use Prolog's DCG grammar. This is because we found it too constraining (could be because of our lack of experience).

Even though, we tried using it. For example, here is how we tried doing the iteration:

```
1 s(N) --> a(N).
2 a(0) --> [].
3 a(R) --> [R], [*], [a].
4 a(M) --> [a], a(N), {M is N + 1}.
5
6 eval(X) :-
7 s(_, X, []).
```

This works for this pattern : $kkk...k \rightarrow N*(k)$, where k is only a **single** character.

We couldn't get it to work with strings or other types (atoms, etc). Even, when using it in combination with phrase/2 like so:

```
1 eval(X, R * (K)) :- phrase(s(R, K), X).
```

LKR Page 14 on 17

5.2 Alternative to our current iteration algorithm

We tried using only Prolog built-ins to achieve it and we came up with this:

```
iteration(A, N, L) :-
2
     length(A,U),
3
     between (1, U, N),
4
     length(L,N),
5
     maplist(=(\_),L),
6
     append (L, A).
7
   iterations (A, N, L, Last) :-
9
     findall(L, iteration(A, N, L), S),
10
     append(_,[Last], S).
```

This works perfectly fine when presented with an array that contains **only** an iteration pattern. It will fail otherwise.

LKR Page 15 on 17

6 Program Usage

To use our program, you should call the eval/2, with an input list.

Examples of usage:

LKR Page 16 on 17

7 Shortlog of the teamwork

Here's the Github link.

Here's our git shortlog:

```
Axel Mendoza (13):
2
         compare suffix
3
         almost finish iteration, drop not implemented
         Merge branch 'master' of github.com:cpcdoy/
4
            LKR_project
5
         almost finish iteration
6
         Merge branch 'master' of github.com:cpcdoy/
            LKR_project
7
         finish iteration
         iteration + * pattern
9
         fix iteration + problem
10
         fix reverse pattern cut string
11
         Merge branch 'master' of github.com:cpcdoy/
            LKR_project
12
         fix wrong matching pattern
13
         fix reverse on sym
14
         sit done with iteration and symmetry
15
   Chady Dimachkie (15):
16
17
         Add: iterations function (ONLY works when the
            input list contains repetitions)
18
         Add: symmetry (no indices)
19
         Add: label_start +
                                 Fix: sublist does not
            return empty arrays anymore
```

LKR Page 17 on 17

```
Merge branch 'master' of https://github.com/
20
            cpcdoy/LKR_project
21
         Add: indices computation
         Merge branch 'master' of https://github.com/
22
            cpcdoy/LKR_project
         Add: symmetry string output
23
         Add: display symmetry +
24
                                    Add: get longest
            symmetry
         Add: correctly displayed symmetry
25
26
         Clean: function names
27
         Add: SIT.pl which merges everything
28
         Fix: cut_string_sym typo
29
         Fix: get_symmetry returns the original array
            when nothing was found
30
         Add: README.md
31
         Add: grammar.pl
```