

Projet

Programmer par des exemples

Concat est un petit langage de programmation qui permet d'accomplir des tâches simples de transformations de chaînes de caractères. L'objectif de ce projet est de faire de la *programmation par l'exemple* : Au lieu de faire écrire un programme Concat par un programmeur, un utilisateur va donner quelques exemples d'une entrée et de la sortie attendue du programme. Puis, le programme que vous allez écrire va produire un programme Concat qui est cohérent avec les exemples fournis (ou une erreur quand un tel programme n'existe pas).

Puisque votre tâche sera d'engendrer des programmes Concat à partir d'exemples, il est inutile de réfléchir à une syntaxe concrète. Le programme que vous écrirez va construire le programme Concat directement sous la forme d'une syntaxe abstraite, et puis l'afficher sous une forme que vous pouvez choisir. En particulier il est inutile de créer un parseur pour ce langage.

1 Notations

Si s est une chaîne de caractères de longueur n , et $0 \leq i < j \leq n$, alors nous écrivons $s[i..j[$ pour la sous-chaîne de s qui va de l'index i *inclus* à l'index j *exclus*. Par exemple, `"hello"[2..4[= "ll"`. On écrit $s[..j[$ pour $s[0..j[$, et $s[i..[$ pour $s[i..n[$. On remarque que $s[i..j[\cdot s[j..k[= s[i..k[$.

Quand on peut découper une chaîne s en $s = s_1 \cdot s_2 \cdot s_3$, où chacun des s_1, s_2, s_3 peut être vide, on appelle s_1 un préfixe, s_2 un facteur, et s_3 un suffixe de s .

2 Étape 1

Un programme Concat décrit une fonction qui envoie une chaîne de caractères vers une chaîne de caractères. La syntaxe abstraite est donnée à la figure 1, et la sémantique est résumée à la figure 2.

```
program  → expression  ...  expression
expression  → const(string) | extract(pos_expression, pos_expression)
pos_expression  → forward(nat) | backward(nat)
```

Figure 1: Syntaxe abstraite de Concat - Étape 1

$$\begin{aligned}
\llbracket e_1 \cdots e_n \rrbracket(s) &= \llbracket e_1 \rrbracket(s) \cdot \dots \cdot \llbracket e_n \rrbracket(s) \\
\llbracket \text{const}(sc) \rrbracket(s) &= sc \\
\llbracket \text{extract}(p_1, p_2) \rrbracket(s) &= s[\llbracket p_1 \rrbracket(s) .. \llbracket p_2 \rrbracket(s) [\\
\llbracket \text{forward}(i) \rrbracket(s) &= i \\
\llbracket \text{backward}(i) \rrbracket(s) &= \text{length}(s) - i
\end{aligned}$$

Figure 2: Sémantique de Concat - Étape 1

- Un programme est une séquence d'expressions, la sémantique d'un programme $e_1 \dots e_n$ appliqué à une entrée s est la concaténation des sémantiques de chaque expression e_i appliquée à s .
- Une expression est soit const appliqué à une chaîne sc , soit extract de deux pos.expressions p_1 et p_2 . Dans le premier cas la sémantique est la chaîne sc . Dans le deuxième cas, quand la sémantique de p_1 appliquée à s est i_1 et la sémantique de p_2 appliquée à s est i_2 , alors la sémantique de l'expression entière appliquée à s est $s[i_1..i_2[$.
- Une pos.expression est une expression qui dénote une position dans une chaîne de caractères. À l'étape 1 c'est soit forward(i), qui dénote l'index i , ou bien backward(i) qui dénote l'index $n - i$ lors d'une application à une chaîne de longueur n .

Par exemple :

- le programme extract(forward(0), backward(0)) est la transformation identique qui transforme une chaîne s en s .
- le programme suivant:

const("Hello, "); extract(forward(3), backward(7))

transforme la chaîne "Mr Smith junior" en "Hello, Smith".

Le programme que vous devrez écrire aura la fonctionnalité suivante : il lit d'abord dans un fichier une séquence de paires associant à une chaîne d'entrée une chaîne de résultat. Chaque paire constitue une ligne de votre fichier, et entrée et résultat seront séparés par le symbole tabulation. Votre programme doit ensuite calculer un programme Concat p tel que pour toute paire exemple (in, out) , on ait $\llbracket p \rrbracket(in) = out$. Par exemple, étant donné les exemples suivants :

Entrée	Résultat
10/10/2017	10
05-15-2015	15
20.02.2020	02

vosre programme peut trouver le programme Concat suivant :

extract(forward(3), forward(5))

Sur cet exemple, on voit déjà qu'il peut y avoir plusieurs programmes Concat possibles pour un même problème, par exemple en utilisant ici des backward à la place des forward.

L'idée du programme est la suivante : pour chacune des k paires (in_i, out_i) on construit l'ensemble $P_i = \{p_i \mid \llbracket p_i \rrbracket(in_i) = out_i\}$, puis on calcule l'intersection $\bigcap_{i=1}^k P_i$, et finalement on choisit un programme dans cette intersection. Si l'intersection est vide on peut signaler une erreur. Dans l'exemple précédent, l'ensemble P_1 correspondant à la première paire va contenir, parmi d'autres, les programmes suivants :

- `const("10")`
- `extract(forward(3), forward(5))`
- `extract(forward(3), forward(4)); const("0")`

et beaucoup d'autres encore. On se rend rapidement compte du fait qu'il ne va pas être possible d'énumérer *explicitement* tous les éléments d'un tel ensemble P_i , car si la chaîne résultat a une longueur n alors il y a déjà 2^n possibilités de la découper en morceaux, et pour chaque morceau du découpage il y aura plusieurs possibilités de l'obtenir soit comme constante, soit par l'extraction de l'entrée à des positions différentes.

Pour cette raison, au lieu d'énumérer explicitement tous les éléments d'un ensemble P_i , on va plutôt en construire une représentation compacte. Cette représentation doit permettre les deux opérations suivantes :

- calculer l'intersection de deux ensembles, nécessaire pour combiner les ensembles obtenus pour des paires différentes ;
- extraire un programme concret, si possible le programme le plus simple qui se trouve dans un ensemble.

2.1 La représentation d'un ensemble de programmes

Nous allons utiliser ici une représentation par des DAG (*directed acyclic graphs* - graphes orientés non cycliques) afin de représenter tous les ensembles de programmes P_i vus ci-dessous, et de représenter leurs intersections. Notez que cette représentation est spécifique à notre problème: elle ne permettra pas d'encoder un ensemble arbitraire de programmes. On décrit d'abord le DAG qui représente l'ensemble P qui correspond à une chaîne d'entrée in de longueur n et un résultat out de longueur m .

L'ensemble de nœuds du graphe est $\{0, \dots, m\}$, et il y a une arête de i vers j quand $i < j$. Le nœud 0 est la source, et m est le puits de ce graphe (voir la figure 3).

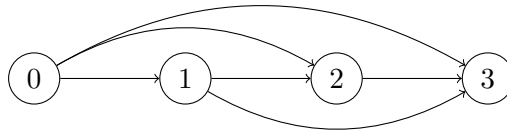


Figure 3: Le graphe des positions pour $m = 3$, ici montré sans les étiquettes

Une arête de i vers j va être étiquetée par une représentation de l'ensemble des expressions e (au sens de la figure 1) tel que $\llbracket e \rrbracket(in) = out[i..j]$. Par conséquent, tout programme p avec $\llbracket p \rrbracket(in) = out$ correspond à un chemin dans le graphe de la source vers le puits, et les expressions qui constituent le programme p sont des expressions qu'on peut trouver sur les arêtes de ce graphe.

L'arête de i vers j est décorée avec un ensemble. Dans cet ensemble on trouvera toujours $\text{const}(\text{out}[i..j])$, les autres éléments sont des extract appliqués à deux *ensembles* de position. Si la chaîne $\text{out}[i..j]$ apparaît plusieurs fois dans la chaîne entrée alors vous construisez un extract pour chaque occurrence de $\text{out}[i..j]$ dans in . Et ce extract aura comme arguments les ensembles de pos_expressions p donnant respectivement les index du début et de la fin de cette occurrence de $\text{out}[i..j]$ dans in . Notez que pour l'instant ces ensembles de positions sont de taille 2 (un forward et un backward) mais cela changera par la suite. En résumé, la décoration de l'arête de i vers j est

$$\{\text{const}(\text{out}[i..j])\} \cup \bigcup_{\text{in}[k..l]=\text{out}[i..j]} \text{extract}(\{p \mid \llbracket p \rrbracket(\text{in}) = k\}, \{p \mid \llbracket p \rrbracket(\text{in}) = l\})$$

Note : dans la syntaxe abstraite de Concat, extract prend deux pos_expressions en arguments, mais dans les étiquettes d'un DAG, extract prend deux *ensembles* de pos_expressions en arguments.

2.1.1 Exemple

Pour rendre les exemples plus lisibles nous écrivons ici \underline{f} et \underline{b} au lieu de forward et backward.

1. $\text{in} = \text{"abad"}$, $\text{out} = \text{"dxa"}$. Le graphe a les 4 nœuds 0, 1, 2, 3 (figure 3), avec les arêtes suivantes :

de	vers	étiquette
0	1	$\text{const}(\text{"d"}), \text{extract}(\{\underline{f}(3), \underline{b}(1)\}, \{\underline{f}(4), \underline{b}(0)\})$
0	2	$\text{const}(\text{"dx"})$
0	3	$\text{const}(\text{"dxa"})$
1	2	$\text{const}(\text{"x"})$
1	3	$\text{const}(\text{"xa"})$
2	3	$\text{const}(\text{"a"}), \text{extract}(\{\underline{f}(0), \underline{b}(4)\}, \{\underline{f}(1), \underline{b}(3)\}), \text{extract}(\{\underline{f}(2), \underline{b}(2)\}, \{\underline{f}(3), \underline{b}(1)\})$

2. $\text{in} = \text{"efegh"}$, $\text{out} = \text{"ghxe"}$ Le graphe a les 5 nœuds 0, 1, 2, 3, 4, avec les arêtes suivantes :

de	vers	étiquette
0	1	$\text{const}(\text{"g"}), \text{extract}(\{\underline{f}(3), \underline{b}(2)\}, \{\underline{f}(4), \underline{b}(1)\})$
0	2	$\text{const}(\text{"gh"}), \text{extract}(\{\underline{f}(3), \underline{b}(2)\}, \{\underline{f}(5), \underline{b}(0)\})$
0	3	$\text{const}(\text{"ghx"})$
0	4	$\text{const}(\text{"ghxe"})$
1	2	$\text{const}(\text{"h"}), \text{extract}(\{\underline{f}(4), \underline{b}(1)\}, \{\underline{f}(5), \underline{b}(0)\})$
1	3	$\text{const}(\text{"hx"})$
1	4	$\text{const}(\text{"hxe"})$
2	3	$\text{const}(\text{"x"})$
2	4	$\text{const}(\text{"xe"})$
3	4	$\text{const}(\text{"e"}), \text{extract}(\{\underline{f}(0), \underline{b}(5)\}, \{\underline{f}(1), \underline{b}(4)\}), \text{extract}(\{\underline{f}(2), \underline{b}(3)\}, \{\underline{f}(3), \underline{b}(2)\})$

2.2 Intersection de deux ensembles

Considérons un DAG D_1 avec un ensemble de nœuds N_1 , source s_1 et puits t_1 , et un DAG D_2 avec ensemble de nœuds N_2 , source s_2 et puits t_2 . Leur intersection est définie ainsi :

- L'ensemble de nœuds est $N_1 \times N_2 = \{(n_1, n_2) \mid n_1 \in N_1, n_2 \in N_2\}$.

2. Nous continuons avec les nœuds utiles trouvés à l'étape précédente :

de	vers	étiquette
(1, 1)	(2, 2)	<i>vide</i>
(1, 1)	(2, 3)	<i>vide</i>
(1, 1)	(3, 4)	<i>vide</i>
(1, 2)	(2, 3)	<u>const</u> ("x")
(1, 2)	(3, 4)	<i>vide</i>

3. Il nous reste un seul nœud utile pour continuer :

de	vers	étiquette
(2, 3)	(3, 4)	<u>extract</u> ({f(0)}, {f(1)}), <u>extract</u> ({f(2)}, {f(3)})

2.3 Extraire un programme d'un ensemble

Extraire un programme d'un DAG est maintenant trivial. Pour chaque arête, décorée avec un ensemble de const et de extract, vous gardez seulement l'expression la plus simple selon un critère qui vous semble pertinent. Puis, pour chaque arête vous calculez un poids basé sur l'unique décoration, de sorte que ce poids soit plus important pour des décorations plus complexes. Encore une fois, vous pouvez choisir librement votre fonction de poids. Ce choix deviendra par contre plus important à l'étape 2. Enfin, vous extrayez de ce graphe un plus court chemin de la source au puits selon l'algorithme bien connu de Dijkstra. Le programme Concat choisi est alors la séquence des expressions sur les arêtes de ce chemin.

2.3.1 Exemple (suite de 2.2.1)

Dans notre exemple, il y a un seul chemin de la source au puits : $(0, 0) \rightarrow (1, 2) \rightarrow (2, 3) \rightarrow (3, 4)$. Un programme Concat possible est donc

extract(forward(3), backward(0)); const("x"); extract(forward(0), forward(1))

Il y a une autre possibilité :

extract(forward(3), backward(0)); const("x"); extract(forward(2), forward(3))

3 Étape 2

La seule chose qui change à l'étape 2 est la définition des `pos_expressions`: on doit maintenant pouvoir spécifier également des positions dans l'entrée par une forme restreinte d'expressions régulières.

```

pos_expression → forward(nat) | backward(nat)
                | before(regexp) | after(regexp)
                | beforelast(regexp) | afterlast(regexp)
regexp → token...token
token → plus(class) | neg(class)
class → alphanumeric | numeric | alpha | lower | upper | special(char)

```

Figure 4: Extension de la syntaxe abstraite à l'étape 2

3.1 Les classes de caractères

Chaque **class** décrit une classe de caractères : les alpha-numériques, les chiffres, les lettres de l'alphabet, les lettres minuscules et les lettres majuscules. Une classe de la forme special(*c*), où le caractère *c* est un caractère spécial, est une classe constituée de cet unique caractère *c*. Vous pouvez choisir quel ensemble de caractères spéciaux vous traitez, par exemple $\{/, -, (,)\}$, mais attention à l'explosion combinatoire si cet ensemble est trop grand.

3.2 Les tokens

Soit *cl* une classe. Le token plus(*cl*) filtre un mot *m* quand *m* est non-vidé, et tous les caractères de *m* sont dans *cl*. Le token neg(*cl*) filtre un mot *m* quand *m* est non-vidé, et *aucun* caractère de *m* n'est dans *cl*.

3.3 Les expressions régulières

Une expression régulière est une séquence non-vidé de tokens. L'expression régulière vide (i.e. constituée de 0 token) ne sera pas acceptée dans nos programmes. Le filtrage d'un mot par une expression régulière est défini par l'algorithme glouton suivant :

- Une expression régulière constituée d'un unique token *t* filtre un mot *m* lorsque *t* filtre *m*.
- Une expression régulière constituée de plusieurs tokens peut s'écrire *tr*, avec *t* un token et *r* une expression régulière (non-vidé). Alors cette expression régulière *tr* filtre un mot *m* si on peut décomposer *m* en $m = m_1 \cdot m_2$ tels que :
 - *t* filtre m_1
 - *t* ne filtre aucun préfixe de m_2
 - *r* filtre m_2

Vu la définition, lorsqu'une expression régulière filtre un mot, ce mot est forcément non-vidé.

3.4 Extension des pos_expressions

Commençons par décrire before(*r*) et after(*r*). Étant donnée une expression régulière *r* et une chaîne de caractères *s*, soit *i* un index minimal de *s* tel que *r* filtre un préfixe de $s[i..]$, et pour ce *i* soit *j* maximal tel que *r* filtre $s[i..j]$. En d'autres mots, on cherche le premier facteur de *s* qui est

filtré par r , et ce facteur est prolongé autant que possible. La pos_expression before(r), quand appliquée à cette chaîne s , dénote la position i , et la pos_expression after(r), quand appliquée à s , dénote la position j .

Pour beforelast(r) et afterlast(r), on reprend la même idée, mais cette fois on cherche le *dernier* facteur de s filtré par r . Plus précisément, soit j un index maximal de s tel que r filtre un suffixe de $s[..j]$, et pour ce j soit i minimal tel que r filtre $[i..j]$. La pos_expression beforelast(r), quand appliquée à cette chaîne s , dénote la position i , et la pos_expression afterlast(r), quand appliquée à s , dénote la position j . Voir l'exemple fourni `remove-last-dir` pour une situation qui nécessite un beforelast.

3.5 Conseils

Dans votre programme vous auriez besoin de connaître, pour tous i et j , l'ensemble des expressions régulières qui filtrent $s[i..j]$. Vous calculerez cette information de façon ascendante (bottom-up) en vous inspirant de la programmation dynamique.

Vous pouvez gagner en efficacité en déterminant d'avance un ensemble pertinent de classes de caractères en fonction du problème à traiter. Par exemple, une classe qui contient tous les caractères d'entrée, ou qui n'en contient aucun, ne nous sert à rien. De plus, deux classes de caractères cl_1 et cl_2 seront redondantes pour notre problème lorsque pour tout caractère c de l'entrée, $c \in cl_1$ ssi $c \in cl_2$. Il suffit alors de considérer une seule des deux classes. Comme on dispose de tokens négatifs, la situation $c \in cl_1$ ssi $c \notin cl_2$ pour tout caractère c de l'entrée mènent également à une redondance de classes.

4 Pour aller encore plus loin

L'idée de ce projet est tiré de l'article [Gul11]. Vous y trouverez des idées d'extensions possibles. Une première extension possible est de considérer des pos_expressions de la forme after(**regex**, n) ou before(**regex**, n), où n est un nombre naturel, qui dénotent le n -ième facteur du mot filtré par l'expression régulière. Vous trouverez dans cet article en plus des autres extensions comme par exemple des programmes avec des boucles, ou des distinctions de cas.

5 Organisation

Voir le fichier `projet/README.md` dans le dépôt git du cours.

References

- [Gul11] Sumit Gulwani. Automating string processing in spreadsheets using input-output examples. In Thomas Ball and Mooly Sagiv, editors, *Proceedings of the 38th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL)*, pages 317–330, Austin, TX, USA, January 2011. ACM.