

Natural to Python Compiler

Nino Jeannet et Sébastien Peiris

Présenté à Fabrizio Albertetti

12 Janvier 2020

Contents

1	Introduction	2
1.1	Présentation	2
1.2	Langage Natural	2
1.3	Natural to Python Compiler	2
2	Analyse lexicale	2
2.1	Mots réservés	2
2.2	Lexèmes	3
2.3	Automates d'états finis	3
2.3.1	Variable	3
2.3.2	Assignment	3
2.3.3	Assignment et calcul	4
3	Analyse syntaxique	4
3.1	Définition de variable	4
3.2	Instanciation de variable	4
3.3	Boucle FOR	4
3.4	Bloc IF...ELSE	5
3.5	Affichage	5
3.6	Problèmes de syntaxe	5
3.6.1	Boucle FOR	5
3.6.2	Condition	5
3.6.3	"=" et "=="	5
3.7	Nombres	5
3.7.1	Nombres décimaux	5
3.8	Arbres syntaxiques	6
3.8.1	Boucle FOR	6
3.8.2	Condition	7
4	Grammaire	8
4.1	Exemples de code	8
4.1.1	Instanciation de variables	8
4.1.2	Bloc de condition	9
4.1.3	Exemple complexe	9
5	Analyse sémantique	10
5.1	Couture	10
5.1.1	Exemples de coutures	10
5.2	Gestion d'erreurs	11
6	Compilateur	12
7	Problèmes rencontrés	12
8	Conclusion	12

1 Introduction

1.1 Présentation

Ce projet a été réalisé dans le cadre du cours de Compilateur. Il a pour but la création d'un langage de programmation ainsi que sa compilation vers un langage de haut niveau. Ce compilateur est intégralement réalisé en *Python* et les analyseurs lexical et syntaxique sont basés sur *PLY*.

1.2 Langage Natural

Natural est donc le langage que nous avons créé. Sa principale caractéristique est son très haut niveau. En effet, celui-ci est proche du langage français parlé. Son but serait d'être compréhensible pour des personnes n'ayant pas ou peu de bases en programmation. Étant donné la complexité de la langue française, nous avons été forcés de limiter l'étendue du langage. En effet, implémenter un compilateur qui comprend le langage "parlé" et qui puisse le traduire instantanément en *bytecode* Python relève de l'impossible. Il nous a donc fallu déterminer des règles et des limites à *Natural*.

1.3 Natural to Python Compiler

Notre compilateur a donc pour but de compiler du code *Natural* en python. *Natural* est capable de:

- Déclarer et instancier des variables
- Effectuer des opérations entre variables (*/-/+)
- Effectuer des boucles *for*
- Effectuer des test *if / else*
- Afficher du contenu

2 Analyse lexicale

L'analyseur lexical a pour but de définir les différents lexèmes du langage *Natural* ainsi que d'autres modalités telles que les mots réservés. Ceux-ci ne pourront pas servir de nom de variable. L'analyse lexicale s'effectue dans le fichier *lex.py*.

2.1 Mots réservés

Les mots réservés du langage sont les suivants:

```
egal, a, vaut, plus, moins, fois, divise, par, pour, allant, de, pas, afficher,  
si, inferieur, superieur, sinon, debut, fin
```

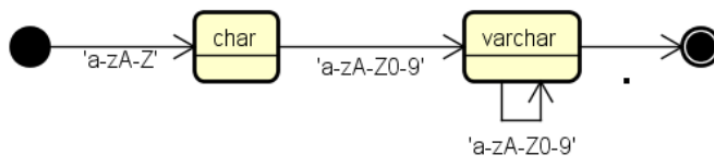
2.2 Lexèmes

Nous avons défini différents lexèmes (tokens) qui composeront notre langage. Ces lexèmes sont définis à l'aide d'expressions régulières. Nous avons défini les lexèmes suivants :

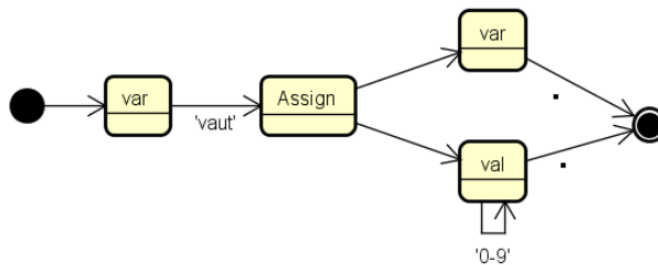
- NUMBER
- ADD_OP
- MUL_OP
- COMPARABLE
- TO
- STEP
- IDENTIFIER
- STRING
- + la liste des mots réservés

2.3 Automates d'états finis

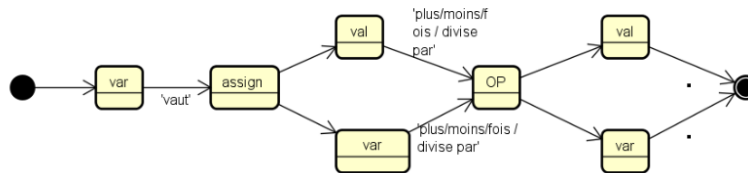
2.3.1 Variable



2.3.2 Assignment



2.3.3 Assignment et calcul



3 Analyse syntaxique

Lors de la création de notre langage, nous avons dû définir une syntaxe à respecter pour que l'on arrive à le compiler. Pour ceci, nous avons défini une syntaxe à respecter pour chaque type d'éléments.

3.1 Définition de variable

Un nom de variable ne peut contenir que des chiffres et des lettres et au minimum un caractère. Il doit également débuter par une lettre (majuscule ou minuscule). L'expression régulière correspondant est la suivante: $[A - Z a - z][A - Z a - z 0 - 9]^+$

3.2 Instanciation de variable

Pour déclarer et instancier une variable il faut utiliser la syntaxe :

```
<nom_de_variable> vaut <valeur_de_variable>.
```

On peut instancier une variable, soit avec une valeur numérique, soit avec une expression qui vaut une valeur numérique.

3.3 Boucle FOR

La boucle *FOR* s'écrit de la manière suivante :

```
1 pour <var> allant de <depart> a <arrivee> par pas de <pas>.
2 debut
3     # do something
4 fin.
```

Cette syntaxe va faire une boucle en initialisant la variable d'incrémentation avec la valeur de départ, en l'incrémentant du pas à chaque itération tant qu'il est plus petit que la valeur d'arrivée.

Important: il n'est pas nécessaire d'utiliser une tabulation avant le contenu de la boucle mais il est fortement recommandé de l'utiliser pour des raisons de lisibilité.

3.4 Bloc IF...ELSE

Pour utiliser un bloc de condition(if..else), il faut utiliser la syntaxe suivante :

```
1 si <condition>.  
2 debut  
3     # do something  
4 sinon  
5     #do something else  
6 fin.
```

3.5 Affichage

Pour afficher une valeur, il faut utiliser la syntaxe ci-dessous:
Si on affiche une chaîne de caractère :

```
afficher '<chaîne>'
```

Si on affiche une valeur :

```
afficher <valeur>
```

3.6 Problèmes de syntaxe

Durant la définition des ces différentes syntaxes, nous avons détecté plusieurs cas qui pourraient poser problème dans la suite du projet. Voici quelques exemples de cas où nous avons adaptés notre syntaxe pour éviter toute ambiguïté.

3.6.1 Boucle FOR

Nous avons, au départ, décidé d'utiliser le symbole "-" pour définir les lignes contenues dans la boucle sans avoir besoin de délimiteur. Mais après réflexion, cette méthode ne fonctionne pas dans le cas de boucles imbriquées car nous aurions des enchaînements de "-" et toute lisibilité serait perdue. Nous avons donc choisi d'utiliser des délimiteurs (*debut* et *fin*).

3.6.2 Condition

Idem que pour la boucle for, nous utilisons des délimiteurs plutôt que des "-".

3.6.3 "=" et "=="

Afin de différencier le = du == nous avons simplement décidé d'utiliser des lexème unique pour les différencier. Donc nous avons *egal a* qui correspond à == et *vaut* qui correspond à =.

3.7 Nombres

3.7.1 Nombres décimaux

Afin de d'éviter la confusion avec les fins de lignes, nous avons décidé d'utiliser la virgule dans les nombres décimaux. Exemple :

```
x vaut 13,5
```

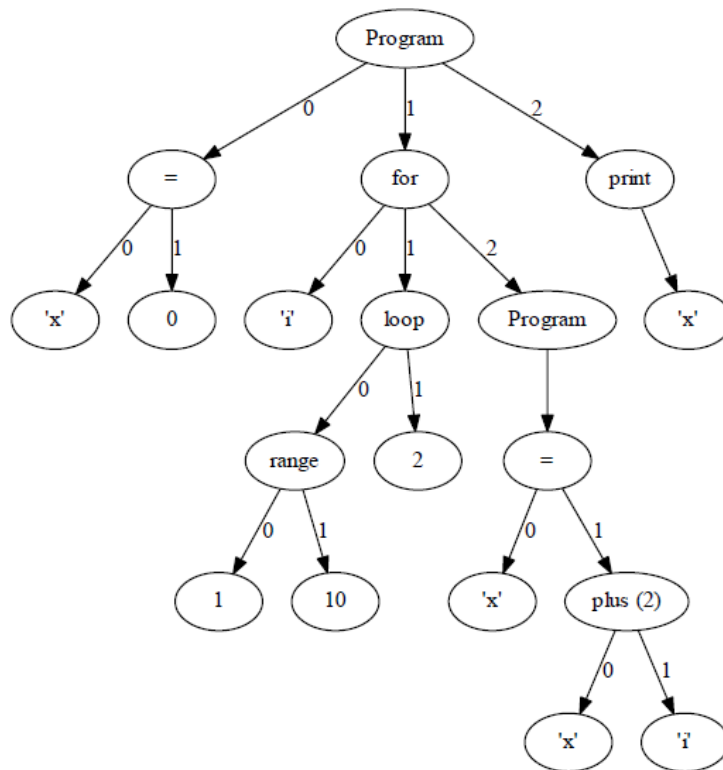
3.8 Arbres syntaxiques

3.8.1 Boucle FOR

Code *Natural*:

```
1 x vaut 0.  
2 pour i allant de 1 a 10 par pas de 2.  
3 debut  
4   x vaut x plus i.  
5 fin.  
6 afficher x.
```

Arbre syntaxique:

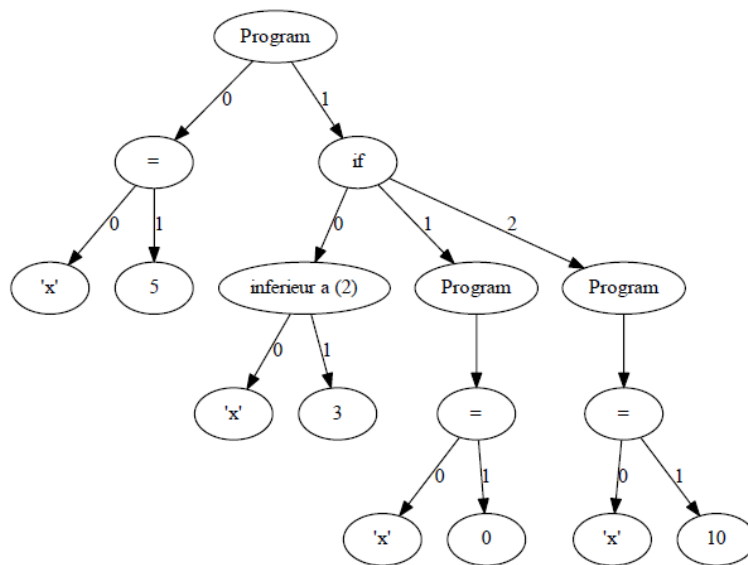


3.8.2 Condition

Code *Natural*:

```
1 x vaut 5.  
2 si x inferieur a 3.  
3 debut  
4   x vaut 0.  
5 sinon  
6   x vaut 10.  
7 fin.
```

Arbre syntaxique:



4 Grammaire

Après avoir défini les différentes syntaxes, on peut maintenant tout regrouper pour définir notre grammaire. Cette grammaire sera notamment utile lors de l'analyse syntaxique afin de construire un arbre syntaxique à partir d'un code source.

```

PROGRAMME -> STATEMENT '.'
PROGRAMME -> STATEMENT '.' PROGRAMME

STATEMENT -> ASSIGNATION | STRUCTURE
STATEMENT -> afficher EXPRESSION | afficher EXPRESSIONSTRING

STRUCTURE -> pour ITERATEUR to LOOP '.' debut PROGRAMME fin
ITERATEUR -> IDENTIFIIFER
LOOP -> RANGE STEP EXPRESSION
RANGE -> EXPRESSION a EXPRESSION

STRUCTURE -> si BOOLEAN '.' debut PROGRAMME sinon PROGRAMME fin
BOOLEAN -> EXPRESSION comparable EXPRESSION

CONDITION -> EXPRESSION inferieur a EXPRESSION
CONDITION -> EXPRESSION superieur a EXPRESSION
CONDITION -> EXPRESSION egal a EXPRESSION

EXPRESSIONSTRING -> STRING
EXPRESSION -> NUMBER | IDENTIFIER
EXPRESSION -> EXPRESSION ADD_OP EXPRESSION
EXPRESSION -> EXPRESSION MUL_OP EXPRESSION
EXPRESSION -> EXPRESSION SUB_OP EXPRESSION
EXPRESSION -> EXPRESSION DIV_OP EXPRESSION

ASSIGNATION -> IDENTIFIER vaut EXPRESSION

```

4.1 Exemples de code

4.1.1 Instanciation de variables

Python :

```

1 x = 0
2 y = x + 12
3
4 for i in range(0, 10, 1):
5     x = 2 * y
6     y = y - 1
7 print(x)
8 print(y)

```

Natural :

```
1 x vaut 0.
2 y vaut x plus 12.
3 pour i allant de 0 a 9 par pas de 1.
4 debut
5     x vaut 2 fois y.
6     y vaut y moins 1.
7 fin.
8 afficher x.
9 afficher y.
```

4.1.2 Bloc de condition**Python :**

```
1 if x < y:
2     print('success')
3 else:
4     print('error')
```

Natural:

```
1 si x inferieur a y.
2 debut
3     afficher 'success'.
4 sinon
5     afficher 'error'.
6 fin.
```

4.1.3 Exemple complexe**Python :**

```
1 x = 1
2 y = 2 + 12
3
4 for i in range(0, 10, 1):
5     if x < y:
6         print(x-y)
7     else:
8         for i in range(1,3,2):
9             x = x + 1
10    print(x)
```

Natural:

```
1 x vaut 1.
2 y vaut 2 plus 12.
3 pour i allant de 0 a 10 par pas de 1.
4 debut
5     si x inferieur a y.
6     debut
7         afficher x moins y.
8     sinon
9         pour i allant de 1 a 3 par pas de 2.
10        debut
11            x vaut x plus 1.
12        fin.
13    fin.
14    afficher x.
15 fin.
```

5 Analyse sémantique

Une partie d'analyse sémantique est nécessaire afin d'effectuer certaines vérifications et empêcher un certain nombre d'erreurs en afficher des messages d'avertissement. Étant donné que le projet se limite à la réalisation d'un compilateur, il n'est pas nécessaire d'arrêter l'exécution si une erreur apparaît. Nous avons donc choisi d'afficher des *warnings* mais de tout de même effectuer l'opération de compilation jusqu'au bout.

5.1 Couture

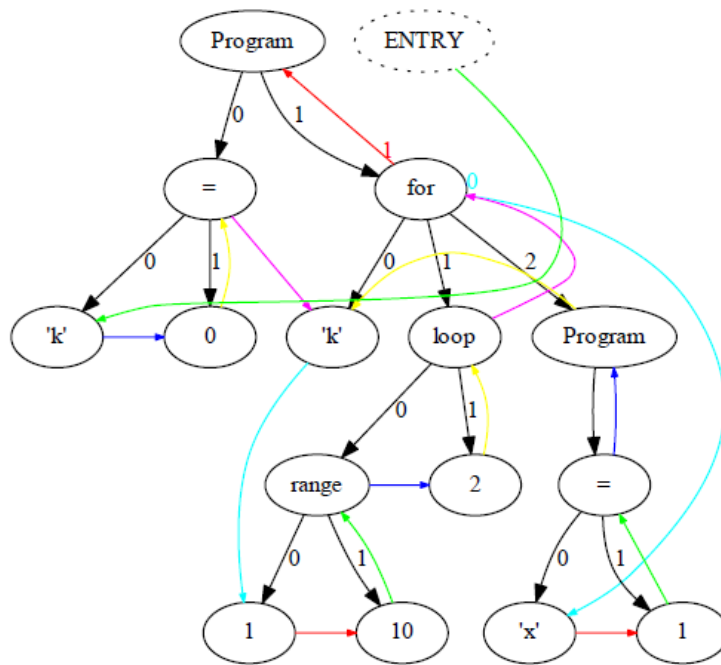
Afin de détecter les potentielles erreurs, il est nécessaire d'implémenter la couture de l'arbre syntaxique. Cette couture parcourt l'arbre syntaxique créé par le parseur et traduit le code noeud par noeud.

5.1.1 Exemples de coutures

Code Natural :

```
1 k vaut 0.
2 pour k allant de 1 a 10 par pas de 2.
3 debut
4     x vaut 1.
5 fin.
```

Couture de l'arbre syntaxique :



5.2 Gestion d'erreurs

Nous avons décidé de gérer 3 types d'erreurs pour lesquels nous avons généré des *warnings* afin d'en informer l'utilisateur :

- L'utilisation de variable non instanciée : Si une variable est utilisée sans être instanciée, on soulève un *warning*. Pour vérifier ceci, nous sauvegardons chaque variable dans un dictionnaire au moment de son instantiation. A chaque token, nous vérifions que la variable existe déjà.
- L'utilisation d'une variable en dehors de son *scope* (pour les boucles for) : Cette fonctionnalité fonctionne avec plusieurs boucles for dans le programme et une imbrication maximum de 2. A chaque fois que nous instancions une variable, la "position" à laquelle elle a été instanciée est sauvegardée. Nous sauvegardons entre autre, l'*id* de la boucle et la profondeur dans laquelle a eu lieu l'instanciation.
- La division par zéro : Pour vérifier la division par zéro, nous vérifions simplement dans un OpNode, que l'opération soit une division et que le second enfant ne soit pas égal à zéro.

6 Compilateur

La partie du compilateur a pour but d'effectuer la traduction finale du code source en *Natural* vers le code de sortie en *Python*. Pour ce faire, le module *compiler.py* parcourt chaque noeud et compile ses enfants afin de traduire le code en langage *Python*.

7 Problèmes rencontrés

Vérification du *scope* pour la condition *if* :

Dans l'analyse sémantique, nous avons voulu vérifier qu'une variable instanciée dans un block *if* soulève un warning si elle est utilisée en dehors de son *scope* (dans le *else* ou block parent). L'implémentation fonctionnait bien avec une simple condition *IF* mais demandait trop de tests avec des conditions imbriquées et plusieurs conditions au même niveaux dans le programme. Il fallait gérer la profondeur, l'*id* et la position(*if* ou *else*) dans le bloc, de la condition. Nous avons donc préféré ne pas laisser cette fonctionnalité incomplètement implémentée dans le programme.

8 Conclusion

Le but final de ce projet était, dans un premier temps, de parvenir à créer un langage de programmation en définissant sa syntaxe et ses règles de grammaire afin de, dans un deuxième temps, le compiler vers un langage de haut niveau existant (ici en *Python*). La démarche suivie a été celle vue en cours, notamment à l'aide des TPs. Nous avons donc suivi les étapes d'analyse lexicale, syntaxique et sémantique avant d'implémenter le compilateur.

Dans la partie dédiée à l'**analyse lexicale**, nous avons défini les différents lexèmes du langage ainsi que les mots réservés. Nous avons également représenté les automates d'états finis pour quelques cas tels que la l'assignation de variables, par exemple.

Lors de l'**analyse syntaxique**, nous avons du définir la syntaxe que chaque élément du code (nombres, blocs *FOR/ IF*, délimiteurs, etc...). Il a également fallu implémenter la construction des arbres syntaxiques, éléments indispensable aux étapes suivantes.

L'**analyse sémantique** consiste en l'implémentation de la couture de l'arbre syntaxique réalisé par l'analyseur syntaxique. Cette couture est essentielle afin de pouvoir parcourir le code et définir s'il contient des incohérences, telles que des variables utilisées mais non déclarées ou encore des divisions pas *zero*.

Au vu des résultats obtenus, nous pouvons constater que l'objectif de la réalisation du compilateur *Natural* vers *Python* a été atteint. Nous obtenons en sortie un fichier script python parfaitement exécutable en ligne de commande et correct syntaxiquement parlant. Ce projet nous a permis de nous rendre compte de la manière dont est implémenté un compilateur, outil que nous utilisons depuis le début de notre formation sans pour autant en comprendre le fonctionnement. Nous avons également pu nous familiariser avec les notions d'analyse lexicale, syntaxique et sémantique, et plus généralement, avec l'analyse de textes en informatique.