

Algorithme Earley 2016/2017

Rapport

*Abderrazak ZIDANE*¹

1. zidane.rezzak@gmail.com

Introduction

Durant mon parcours d'étudiant en informatique, j'ai tout d'abord découvert les compilateurs, puis j'ai appris à les aimer, et aujourd'hui je suis amené à les construire et à les développer. Ce projet est une grande opportunité pour moi d'agrandir mes acquis et connaissances, et de comprendre les aspects théorique et logiciels tout en faisant ce à quoi j'aspire. De plus, les problèmes rencontrés durant ce travail, mon appris à faire de la recherche et à lire de la documentation et thèse.

Depuis le séminaire de Donald Knuth[6] sur l'analyse syntaxique LR en 1960, puis les travaux de DeRemer[1, 2] pour l'extension vers LALR, nous sommes capable de générer automatiquement des analyseurs syntaxiques pour une grande variété de grammaire non contextuelle. Par contre, plusieurs analyseur syntaxique sont écrit manuellement, car souvent, on a pas le luxe de concevoir une grammaire adapter a un générateur d'analyseur syntaxique. Mais aussi, c'est très claire que les concepteurs de langage informatique, n'écrivent pas naturellement des grammaire LR(1).

Une grammaire, non seulement elle définit la syntaxe du langage, mais aussi, c'est le point d'entrées vers la définition de la sémantique, et souvent la grammaire qui facilite la définition de la sémantique n'ai pas LR(1). Ceci est montré pas le développement de la spécification de JAVA. La première édition de cette spécification[4] montre l'effort mis dans la sémantique pour que la grammaire soit LALR(1), par contre dans la 3ème édition de cette spécification[5], la grammaire est (grandement) ambiguë, et ceci montre la difficulté pour faire les transformations adéquates.

Puisque c'est difficile de construire (ou maintenir) des grammaires LR(1) qui garde la sémantique voulu au départ, les développeurs se sont intéressés à d'autres algorithmes comme CYK[9], Earley[3], GLR[8], qui eux ont été développés pour le traitement de langage naturelle à la base (gère l'ambiguïté).

Quand on utilise la grammaire comme point d'entrée pour la définition de la sémantique, on distingue souvent entre **reconnaisseur syntaxique** qui détermine simplement si un mot appartient ou pas à la grammaire, et **analyseur syntaxique** qui retourne la dérivation détaillée d'un mot si elle existe.

Dans leurs versions de base, l'algorithme CYK et Earley sont des reconnaisseurs syntaxiques, alors que GLR est un analyseur syntaxique. Sauf que l'analyseur syntaxique GLR de Tommita a une complexité polynomiale infinie.

Par contre Elizabeth Scott[7], a créé deux algorithmes d'analyse syntaxique basé sur Earley, ayant une complexité cubique dans le pire des cas.

Nous allons tout d'abord comprendre les méthodes d'Elizabeth Scott, et proposer une application écrite en C++ qui implémente ces méthodes.

Introduction

Qu'est ce que l'algorithme Earley?

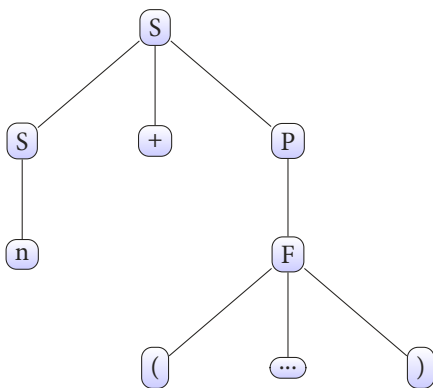
L'algorithme Earley¹, est un parmi les nombreux algorithmes d'analyse syntaxique, est comme tout ces algorithmes, **Earley** a besoin d'une grammaire

$$\begin{array}{lcl} S & = & S + P \\ | & & P \\ P & = & P * F \\ | & & F \\ F & = & (S) \\ | & & n \end{array}$$

... pour transformer une chaine de mot ...

n	+	(n	*	n	+	n)
---	---	---	---	---	---	---	---	---

... en un jolie arbre syntaxique (AST)



Jusqu'à présent rien de spéciale comparé au autre algorithme qu'on conné.

Pourquoi devrions-nous nous soucier ?

Le plus grand avantage d'Erley, est sans doute son accessibilité. La plus part des autres algorithmes offre une restriction sur le type de grammaires, utilisé une grammaire récursif gauche et en rentrera dans une boucle infini, utilisé un autre type et l'algorithme ne marchera plus. Bien sûr il y a des contournements qu'on peut faire, mais souvent ça complique d'avantage l'algorithme et rendra le travail plus complexe.

Pour dire simple Earley marche avec tout.

D'une autre part, pour avoir cette généralité, nous devons sacrifier la vitesse. On ne pourra donc pas rivaliser avec des algorithmes comme Flex/Bison en terme de rapidité brute. Ce n'est si grave, puisque

1. text

- Earley a une complexité cubique $O(3)$, dans les pires des cas, ces même cas qui ne pourront être traités par d'autre algorithme
- La plupart des grammaires simples auront une complexité linéaire
- Même la pire grammaire non-ambigue pourra être analysée en $O(3)$

Vocabulaire

Du point de vue de l'algorithme Earley, la grammaire est constituée de règles. Voici un exemple de règle

$S = S + P$

S est un symbole non terminale, et tout ce qui ne commence pas avec une lettre majuscule est considéré comme étant un symbole terminale (le symbole $+$ dans notre exemple).

Dans le jargon d'Earley, il y a la notion d'items, Voici un item :

$A = B \cdot * C \quad (4)$

C'est juste une règle de grammaire avec des informations en plus, qui représente une reconnaissance partielle.

- Le point représente la position courante qui indique jusqu'où on a parvenu à parser
- Le chiffre 4 représente la position initiale sur l'entrée qu'on veut parser

De plus Earley introduit la notion d'ensemble d'items, chaque ensemble est caractérisé par le fait que les items qui y sont associés, ont la même position courante

Et tous ces ensembles là, sont souvent nommés table Earley

Prédiction, Lecture et Complétion

Pour construire la table Earley, on a besoin de définir trois opérations élémentaires qui s'appliquent sur un item pour produire un autre item :

- Prédiction : le symbole à droite du point est un nom terminal, on ajoute les règles de ce symbole au même ensemble
- Lecture : le symbole à droite du point est un terminal, on regarde si ce symbole coïncide avec la position courante, si oui, on ajoute cet item à l'ensemble suivant.
- Complétion : il n'y a rien à droite du point, et dans ce cas il y a reconnaissance partielle, on regarde l'item parent, et on l'ajoute à cet ensemble

Exemple de construction de table d'Earley

Reprenons cette grammaire :

$S = S + P$
 $| P$
 $P = P * F$
 $| F$
 $F = (S)$
 $| n$

on veut reconnaître l'entrée

n	+	(n	*	n)
---	---	---	---	---	---	---

À l'étape 0, le calcul démarre avec l'ensemble $E(0)$ et les règles de l'axiome ' S '

$E(0)$
$S = \bullet S + P \quad (0)$
$S = \bullet P \quad (0)$

la prédiction du premier item de $E(0)$ nous donnera les mêmes 2 items de $E(0)$, et donc pas besoin de faire quoi que se soit, donc une grammaire récursive gauche ne posera pas de problème à notre algorithme.

La prédiction du deuxième item de $E(0)$ générera deux nouveaux items :

E(0)
S = •S + P (0)
S = •P (0)
P = •P * F (0)
P = •F (0)

Le prédiction du 3ème item de E(0) ne sert a rien. La prédiction du 4ème item de E(0) générera deux nouveaux items supplémentaire :

E(0)
S = •S + P (0)
S = •P (0)
P = •P * F (0)
P = •F (0)
F = •(S) (0)
F = •n (0)

La Lecture du 5ème item de E(0) échoue puisque le symbole ne correspond pas a l'entrée.

La lecture du 6ème item se fait avec succès, est génère un nouveau item dans l'ensemble suivant E(1)

E(1)
F = n• (0)

On a traité tout les items de E(0), attaquons nous a l'ensemble E(1)

La Complétion du premier item de E(1), nous fait ajouter le 4ème item de E(0) a E(1) :

E(1)
F = n• (0)
P = F• (0)

la Complétion du deuxième item de E(1), nous fait ajouter le deuxième et troisième item de E(0) dans E(1)

E(1)
F = n• (0)
P = F• (0)
S = P• (0)
P = P• * F (0)

...

Au finale notre table Earley ressemblera a :

E(0)	E(1)	E(2)	E(3)	E(4)
S = •S + P (0) S = •P (0) P = •P * F (0) P = •F (0) F = •(S) (0) F = •n (0)	F = n• (0) P = F• (0) S = P• (0) P = P• * F (0) S = S• + P (0)	S = S + •P (0) P = •P * F (2) P = •F (2) F = •(S) (2) F = •n (2)	F = (•S) (2) S = •S + P (3) S = •P (3) P = •P * F (3) P = •F (3) F = •(S) (3) F = •n (3)	F = n• (3) P = F• (3) S = P• (3) P = P• * F (3) S = S• + P (3) F = (S•) (2)
E(5)	E(6)	E(7)		
P = P * •F (3) F = •(S) (5) F = •n (5)	F = n• (5) P = P * F• (3) S = P• (3) P = P• * F (3) F = (S•) (2) S = S• + P (3)	F = (S)• (2) P = F• (2) S = S + P• (0)		

On arrive donc a la fin (TODO : condition de reussite)

Que va t'on faire?

Nous allons créer un programme qui va avoir en entrée une grammaire, et aura en sortie un analyseur syntaxique suivant l'algorithme Earley.

Par la suite en va modifier cet algorithme pour notre besoin (TO DO :)

Développement de l'outils

Le programme sera écrit en C++, mais sera facilement traduit dans d'autre langage si nécessaire.

Petite Pré-analyse avant de commencer

Notre outils aura en entrée une grammaire et en sortie, on aura un analyseur syntaxique. Plusieurs questions se sont posées durant le développement de cet outil. Voici un récapitulatif des décisions prises :

- Nous appellerons notre programme `earley`
- La première entrée de notre programme, sera un fichier, qui contiendra la description de la grammaire en format Yacc
- La deuxième entrée de notre programme sera un deuxième fichier, qui contiendra la chaîne à analyser.
- La sortie sera un troisième fichier contenant l'AST
- Les noms de variables, et de fonction seront en anglais, je referai donc au dictionnaire (à faire)

Le programme sera exécuté en ligne de commande suivant la syntaxe suivante :

```
earley <file1> <file2> <file3>
```

`file1` : le fichier de grammaire `file2` : le fichier contenant la chaîne à analyser `file3` : le fichier contenant l'AST

Le Format Yacc

Pour la grammaire que nous fournissons au programme, nous utiliserons le format Yacc simplifié suivant :

Une règle de grammaire a la forme :

`A : BODY ;`

`A` représente un symbole non-terminal, et `BODY` représente une séquence de zéro ou plus de terminaux et non-terminaux. Le colon et le point-virgule sont la ponctuation de Yacc.

un symbole terminal doit être déclaré comme tel au début du fichier :

```
%token n1 n2 p
```

Si un symbole non terminal correspond à la chaîne vide, cela peut être indiqué de manière évidente :

`A ;`

le symbole de départ est considéré comme le côté gauche de la première règle de grammaire dans la section des règles

Bibliographie

- [1] Frank L DeRemer and Thomas J. Pennello. Efficient computation of lalr(1) look-ahead sets. In *ACM Trans. Program. Lang. Syst.*, 4(4) :615–649, October 1982, 1982.
- [2] Franklin L DeRemer. *Practical translators for LR(k) languages*. PhD thesis, Massachusetts Institute of Technology, 1969.
- [3] J Earley. An efficient context-free parsing algorithm. In *Communications of the ACM*, 13(2) :94–102, February 1970.
- [4] Bill Joy James Gosling and Guy Steele. The java language specification. Technical report, Addison-Wesley, 1996.
- [5] Guy Steele James Gosling, Bill Joy and Gilad Bracha. The java language specification third edition. Technical report, Addison-Wesley, 2005.
- [6] Donald E Knuth. On the translation of languages from left to right. In *Information and Control* 8, (6) :607–639, 1965.
- [7] Elizabeth Scott. *SPPF-Style Parsing From Earley Recognisers*. PhD thesis, University of London, 2008.
- [8] Masaru Tomita. *Efficient parsing for natural language*. Kluwer Academic Publishers, Boston, 1986.
- [9] D H Younger. Recognition of context-free languages in time n^3 . In *Inform. Control*, 10(2) :189–208, February 1967.