

Analyse Syntaxique – CTD - Partie 4

Analyse syntaxique ascendante

L'analyse ascendante correspond à la seconde manière de construire systématiquement un automate à pile partant d'une grammaire donnée (voir Cours-TD 4 page 3). Dans cet automate, les transitions empilent les symboles terminaux de la grammaire, et dépilent les non terminaux en ce sens qu'un membre droit de règle situé en sommet de pile est remplacé par le membre gauche correspondant. Plus prosaïquement :

- on part du mot à analyser,
- on remplace itérativement des fragments du mot courant qui correspondent à des membres droits de règle de production par le membre gauche de cette règle,
- l'analyse réussit si le mot courant final est l'axiome de la grammaire

Analyseur ascendant prédictif LR(k)

- **L**eft-to-right scanning : on lit l'entrée de gauche à droite.
- **R**ightmost derivation in reverse : on construit une dérivation droite en partant du mot à analyser.
- on utilise k lettres du mot d'entrée pour faire la prédiction.

Famille d'algorithmes (LR(0), SLR(1), LR(1), LALR(1)) basés sur les opérations de décalage/réduction.

- décalage : on lit des terminaux dans le mot d'entrée
- réduction : on remplace une partie droite par une partie gauche de règle de production.

On souhaite analyser le mot $w = a_1 \dots a_n$ pour une grammaire G . A chaque instant, on dispose :

- d'un suffixe w' du mot d'entrée
- d'une pile contenant des terminaux et des non-terminaux

Initialement, le suffixe est le mot à analyser w et la pile est vide. A chaque étape de l'algorithme, si le suffixe est $a_m \dots a_n$ et le contenu de la pile est $\alpha_1 \dots \alpha_l$ (où α_l est le sommet de pile), on a $\alpha_1 \dots \alpha_l \rightarrow_G^* a_1 \dots a_{m-1}$. L'analyse réussit si le mot w a été entièrement lu et si la pile contient l'axiome.

$S' \rightarrow S$ **Augmentation de la grammaire**

Prenons pour exemple la grammaire suivante : $S \rightarrow Ac$

$A \rightarrow AaAb \mid d$

- Décalage : le premier terminal du mot d'entrée est effacé et est placé au sommet de la pile.

$\boxed{\quad} \quad dadbc \quad \Rightarrow \quad \boxed{d} \quad adbc$

- Réduction : les n premiers symboles sur la pile forme un membre droit de règle de production, on les dépile et on empile le membre gauche correspondant.

$\boxed{d} \quad adbc \quad \Rightarrow \quad \boxed{A} \quad adbc$

En continuant le processus : on ne peut pas réduire, donc on décale

$\boxed{Aa} \quad dbc\$ \quad dadbc \leftarrow Aadbc$

on ne peut pas réduire, donc on décale

$\boxed{Aad} \quad bc\$ \quad dadbc \leftarrow Aadbc$

Ni Aad , ni ad , ni d ne sont préfixe strict d'un membre droit de règle, donc un décalage ne peut être utile ; donc on réduit

$\boxed{AaA} \quad bc\$ \quad dadbc \leftarrow Aadbc \leftarrow AaAbc$

on ne peut pas réduire, donc on décale

$\boxed{AaAb} \quad c\$ \quad dadbc \leftarrow Aadbc \leftarrow AaAbc$

Ni $AaAb$, ni aAb , ni Ab , ni b ne sont préfixe strict d'un membre droit de règle, donc un décalage ne peut être utile ; donc on réduit

$\boxed{A} \quad c\$ \quad dadbc \leftarrow Aadbc \leftarrow AaAbc \leftarrow Ac$

on ne peut pas réduire, donc on décale

$\boxed{Ac} \quad \$ \quad dadbc \leftarrow Aadbc \leftarrow AaAbc \leftarrow Ac$

on ne peut plus décaler ; on réduit

$\boxed{S} \quad \$ \quad dadbc \leftarrow Aadbc \leftarrow AaAbc \leftarrow Ac \leftarrow S$

on réduit à nouveau

$\boxed{S'} \quad \$ \quad dadbc \leftarrow Aadbc \leftarrow AaAbc \leftarrow Ac \leftarrow S \leftarrow S'$

Le mot $dadbc$ est bien engendré par la grammaire.

Comment par un algorithme décider à chaque étape si on doit réduire (et par quelle règle) ou décaler ? On construit un automate déterministe appelé *automate LR(0)* qui dira quoi faire en fonction du contenu de la pile et la première lettre de l'entrée.

Automates LR(0)

Règles pointées Pour toute règle de la grammaire $X \rightarrow \beta\gamma$ (avec $\beta, \gamma \in (\Sigma \cup V)^*$,

$X \rightarrow \beta \bullet \gamma$ est une *règle pointée*.

Une règle pointée $X \rightarrow \beta \bullet \gamma$ est dite *complète* si $\gamma = \epsilon$.

Exemple Pour les règles $S' \rightarrow S, S \rightarrow Ac, A \rightarrow AaAb \mid d$, les règles pointées sont : $S' \rightarrow \bullet S, S' \rightarrow S \bullet, S \rightarrow \bullet Ac, S \rightarrow Ac \bullet,$

$S \rightarrow A \bullet c, A \rightarrow \bullet AaAb, A \rightarrow A \bullet aAb, A \rightarrow Aa \bullet Ab,$

$A \rightarrow AaA \bullet b, A \rightarrow AaAb \bullet, A \rightarrow \bullet d, A \rightarrow d \bullet.$

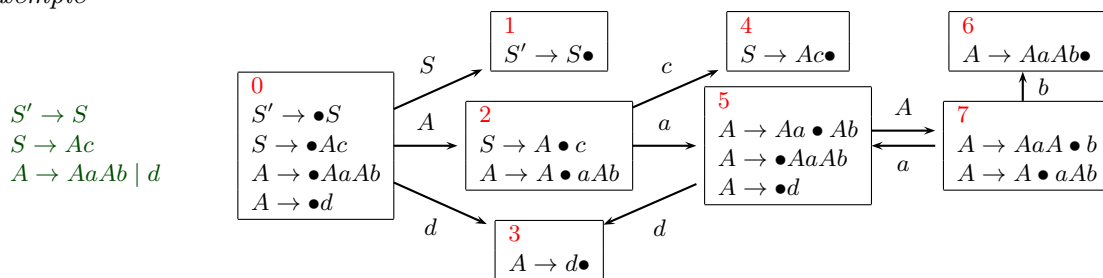
Saturation Un ensemble de règles pointées E est dit *saturé* si $X \rightarrow \beta \bullet Y\gamma$ appartient à E et $Y \in V$ alors pour toute règle $Y \rightarrow \alpha$ de la grammaire, $Y \rightarrow \bullet \alpha$ appartient à E .

Pour un ensemble de règles pointées E , on note $Saturation(E)$ le plus petit ensemble saturé E' contenant E .

Etats de l'automate : Les états Q de l'automate LR(0) seront des ensembles saturés de règles pointées. Ils sont souvent appelés *collections LR(0)* (ou *LR(0)-items*). L' *état initial* de l'automate est $Saturation(\{S' \rightarrow \bullet S\})$.

Transitions de l'automate : Pour un ensemble saturé de règles pointées E et $x \in \Sigma \cup V$, on définit l'état $\delta(E, x)$ comme $Saturation(\{X \rightarrow \beta x \bullet \gamma \mid X \rightarrow \beta \bullet x\gamma \in E\})$ Lors de la construction de l'automate, on ne considèrera bien-sûr que les états accessibles, en partant de l'état initial $Saturation(\{S' \rightarrow \bullet S\})$.

Exemple



La construction de l'automate est totalement indépendante du mot d'entrée (pas de symbole de prévision), d'où son nom : LR(0). On utilise l'automate pour "lire" le mot dans la pile. Si la pile

contient $\alpha_1 \dots \alpha_m$, alors l'état $q = \delta(0, \alpha_1 \dots \alpha_m)$ (où 0 est l'état initial) nous dit ce qu'il faut faire :

- si l'état q contient une règle pointée complète $X \rightarrow \beta\bullet$, alors on réduit en utilisant la règle $X \rightarrow \beta$.
- si le mot courant débute par le terminal a de Σ et s'il existe une transition $\delta(q, a) = q'$, alors on décale.

Si pour tout état q de l'automate,

- q ne contient au plus qu'une seule règle pointée complète (*pas de conflit réduction/réduction*) et
- si q contient une règle pointée complète alors pour aucun a de Σ il n'existe de transition $q' = \delta(q, a)$ (*pas de conflit décalage/réduction*),

alors l'automate est dit LR(0)-consistant et la grammaire est alors LR(0).

Plutôt que d'évaluer systématiquement $\delta(0, \alpha_1 \dots \alpha_m)$ où $\alpha_1 \dots \alpha_m$ est le contenu de la pile (lu du bas vers le haut), on garde trace des états directement dans la pile. La pile est donc une suite $0 \alpha_1 q_1 \alpha_2 \dots q_{m-1} \alpha_m q_m$ où

- les q_i sont des états de l'automate,
- $q_i = \delta(q_{i-1}, \alpha_i) = \delta(0, \alpha_1 \dots \alpha_i)$ pour tout $1 \leq i \leq m$.

Initialement, la pile contient l'état initial 0. Lors d'une étape du calcul,

- si on décale un terminal a , on empile a puis l'état $q = \delta(q_m, a)$.
- si on réduit par une règle $X \rightarrow \alpha_r \dots \alpha_m$, on dépile $q_m, \alpha_m, q_{m-1}, \dots, q_r$ et α_r . on empile X puis l'état $\delta(q_{r-1}, X)$.

Tables LR(0)

L'automate LR(0) est en fait codé dans deux tables, une table *Action* et une table *Successeur* :

- *Action* : $Q \times (\Sigma \cup \{\$\}) \rightarrow$ un ensemble d'actions. Les actions :
 - décaler q : $d q$ avec $q \in Q$.
 - réduire p : $r p$ avec p une production
 - accepter
 - erreur
- *Successeur* : $Q \times V \rightarrow Q$ est la restriction de la fonction de transition aux non-terminaux.

On remplit ces tables à partir de l'automate LR(0).

Table *Action* :

- pour tout $a \in \Sigma, q \in Q$, si $\delta(q, a) = q'$ alors mettre "***d q'***" dans la case (q, a)
- pour tout $a \in (\Sigma \cup \{\$\})$, $q \in Q$, si q contient une règle pointée complète $X \rightarrow \beta\bullet$ avec $X \neq S'$, alors mettre "***r X \rightarrow \beta***" dans la case (q, a)
- mettre "***accepter***" dans la case $(q, \$)$ où q est l'état contenant $S' \rightarrow S\bullet$.
- mettre "***erreur***" dans toutes les cases encore vide.

Si l'automate est LR(0)-consistant alors les cases de la table *Action* contiennent exactement une action.

Table *Successeur* : pour tout $q \in Q$ et $X \in V$, si $\delta(q, X) = q'$ alors mettre q' dans la case (q, X) .

Exemple

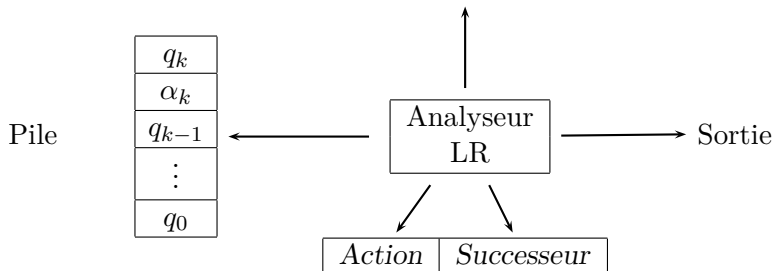
	a	b	c	d	\$
0	erreur	erreur	erreur	d3	erreur
1	erreur	erreur	erreur	erreur	accepter
2	d5	erreur	d4	erreur	erreur
3	r A → d	r A → d	r A → d	r A → d	r A → d
4	r S → Ac	r S → Ac	r S → Ac	r S → Ac	r S → Ac
5	erreur	erreur	erreur	d3	erreur
6	r A → AaAb	r A → AaAb	r A → AaAb	r A → AaAb	r A → AaAb
7	d5	d6	erreur	erreur	erreur

	S'	S	A
0		1	2
1			
2			
3			
4			
5			7
6			
7			

Fonctionnement d'un analyseur LR

Entrée à analyser :

$a_m a_{m+1} \dots a_n$



Ce schéma est le même pour tous les analyseurs LR (LR(0),SLR(1),LALR(1), LR(1),LR(k)) ; seules changent les tables Action et Successeur.

La pile contient

$$0 \alpha_1 q_1 \alpha_2 \dots q_{m-1} \alpha_m q_m$$

et le mot à analyser est $\alpha_{m+1} \dots \alpha_n$.

Si la case (q_m, α_{m+1}) de la table Action contient

- **accepter** alors l'analyse réussit,
- **erreur** alors l'analyse échoue,
- $d q'$, on empile α_{m+1} , puis q' ,
- $r X \rightarrow \alpha_r \dots \alpha_m$, on dépile $\alpha_m, \dots, \alpha_r$ et les états correspondants q_m, \dots, q_r ; on empile X , puis le contenu de la case (q_{r-1}, X) de la table Successeur.

Exemple

0 dabc\$

Action contient d3 en (0, d)

0d3 abc\$

Action contient r A → d en (3, a) et Successeur 2 en (0, A)

0A2 abc\$

Action contient d5 en (2, a)

0A2a5 bc\$

Action contient d3 en (5, d)

0A2a5d3 c\$

Action contient r A → d en (3, b) et Successeur 7 en (5, A)

0A2a5A7 bc\$

Action contient d6 en (7, b)

0A2a5A7b6 c\$

Action contient r A → AaAb en (6, c) et Successeur 2 en (0, A)

0A2 c\$

Action contient d 4 en (2, c)

0A2c4 \$

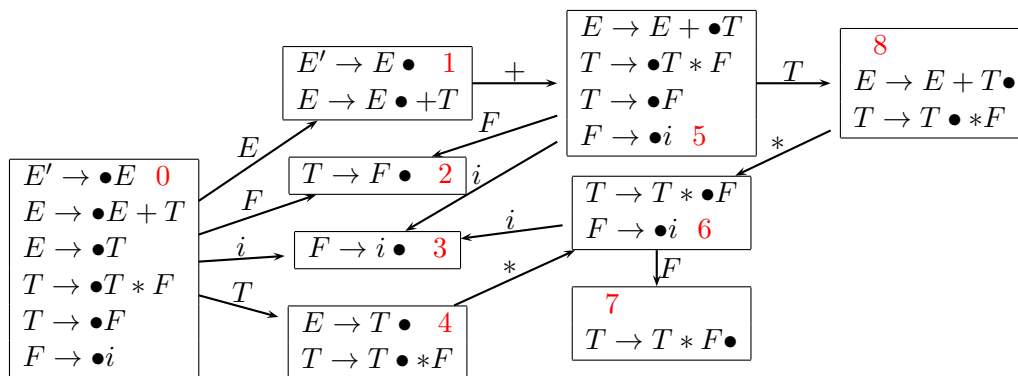
Action contient r S → Ac en (4, \$) et Successeur 1 en (0, S)

0S1 \$

Action contient **accepter** en (1, \$)

Exercice 1. Faire l'automate LR(0) pour la grammaire des expressions suivante :

$$E' \rightarrow E, E \rightarrow E + T \mid T, T \rightarrow T * F \mid T, F \rightarrow i$$



Conflits dans cet automate Cet automate n'est pas LR(0)-consistant : conflit décalage/réduction pour l'état 1 par exemple.

0 $i + i * i \$$ on décale vers l'état 3

0i3 $+ i * i \$$ on réduit par $F \rightarrow i$ et on va dans l'état 2

0F2 $+ i * i \$$ on réduit par $T \rightarrow F$ et on va dans l'état 4

0T4 $+ i * i \$$ on réduit par $E \rightarrow T$ et on va dans l'état 1

0E1 $+ i * i \$$

Conflit entre réduire par $E' \rightarrow E$ et décaler + !! Mais si on réduit alors pour que le mot soit accepté, il faudrait qu'il existe une dérivation (droite) $E' \rightarrow_G^* E' + i * i$; ce n'est bien-sûr pas le cas puisque $Suivant(E') = \{\$ \}$. Donc on décale et on va dans l'état 5.

Quels sont les autres conflits possibles ?

Analyse SLR(1)

Comme on vient de le voir dans l'exercice, on peut résoudre les conflits décalage/réduction en regardant si le terminal en tête du mot restant à analyser est dans l'ensemble des *Suivants* du non-terminal en membre gauche de la règle utilisable pour la réduction. Si c'est le cas, on effectue la réduction, sinon on effectue un décalage. Ceci est formalisé par la méthode SLR(1) dont voici la manière de construire la table *Action* à partir de l'automate LR(0) : Si l'automate LR(0) a une ensemble d'états Q et un ensemble de transitions δ ,

Table Action :

- pour tout $a \in \Sigma, q \in Q$, si $\delta(q, a) = q'$ alors mettre " $d q'$ " dans la case (q, a)
- pour toute règle pointée complète $X \rightarrow \beta \bullet$ de q avec $X \neq S'$, mettre " $r X \rightarrow \beta$ " dans la case (q, a) pour tout $a \in Suivant(X)$.
- mettre "**accepter**" dans la case $(q, \$)$ où q est l'état contenant $S' \rightarrow S \bullet$.
- mettre "**erreur**" dans toutes les cases encore vide.

Table Successeur : identique à la table Successeur LR(0).

Si la table Action contient exactement une action par case, alors la grammaire est SLR(1).

Application à l'exercice précédent

$$\begin{aligned} \text{Suivant}(E') &= \{\$ \} & \text{Suivant}(E) &= \{+, \$\} \\ \text{Suivant}(T) &= \{+, *, \$\} & \text{Suivant}(F) &= \{+, *, \$\} \end{aligned}$$

Table Action

	+	*	i	\$
0	erreur	erreur	d3	erreur
1	d5	erreur	erreur	accepter
2	rT → F	rT → F	erreur	rT → F
3	rF → i	rF → i	erreur	rF → i
4	RE → T	d6	erreur	RE → T
5	erreur	erreur	d3	erreur
6	erreur	erreur	d3	erreur
7	rT → T * F	rT → T * F	erreur	rT → T * F
8	RE → E + T	d6	erreur	RE → E + T

Table Successeur

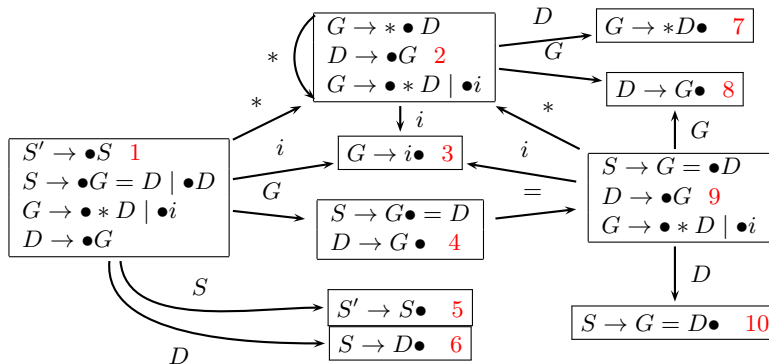
	E'	E	T	F
0		1	4	2
1				
2				
3				
4				
5			8	2
6				7
7				
8				

Une seule action par case dans la table Action ⇒ la grammaire est SLR(1).

- toute grammaire LR(0) est SLR(1)
- toute grammaire SLR(1) est non-ambigüe
- bcp de grammaires non ambigües ne sont pas SLR(1)

Exercice 2. Construire les tables SLR(1) pour la grammaire $\left\{ \begin{array}{l} S \rightarrow G = D \mid D \\ G \rightarrow *D \mid i \\ D \rightarrow G \end{array} \right\}$

Pour cela, on construit l'automate LR(0).



En appliquant la méthode SLR(1) décrite ci-dessus, il va subsister un conflit pour l'état 4 et le terminal = : on peut en effet décaler (d9), mais aussi réduire (rD → G) car le terminal = appartient à Suivant(D) ({=, \$}). Cette grammaire n'est donc pas SLR(1).

Analyse LR(1)

On enrichit les règles pointées avec des terminaux qui permettent de prédire quand effectuer une réduction plutôt qu'un décalage.

Automates LR(1)

Les états d'un automate LR(1) sont des ensembles de règles pointées étendues. Une règle pointée étendue est une "règle pointée associée à un symbole terminal a" ($a \in \Sigma \cup \{\$\}$) et se note $X \rightarrow \alpha \bullet \beta, a$. Si dans un même état, on a $X \rightarrow \alpha \bullet \beta, a$ et $X \rightarrow \alpha \bullet \beta, b$, on notera de manière plus concise $X \rightarrow \alpha \bullet \beta, a \mid b$.

Remarque : si un état contient une règle pointée complète $X \rightarrow \gamma\bullet, a \mid b \mid c$, alors $\{a, b, c\} \subseteq \text{Suivant}(X)$, mais il est possible que ce ne soit qu'un sous-ensemble strict.

Saturation d'un ensemble de règles pointées étendues : un ensemble de règles pointées étendues E est dit *saturé* si $X \rightarrow \beta\bullet Y\gamma, a$ appartient à E et $Y \in V$ alors pour toute règle $Y \rightarrow \alpha$ de la grammaire, $Y \rightarrow \bullet\alpha, b$ avec $b \in \text{Premier}(\gamma a)$ appartient à E .

Pour un ensemble de règles pointées étendues E , on note $\text{Saturation}(E)$ le plus petit ensemble saturé E' contenant E .

Exercice 3. Calculer $\text{Saturation}(\{S' \rightarrow \bullet S, \$\})$ pour la grammaire de l'exercice 2.

Etats de l'automate LR(1) :

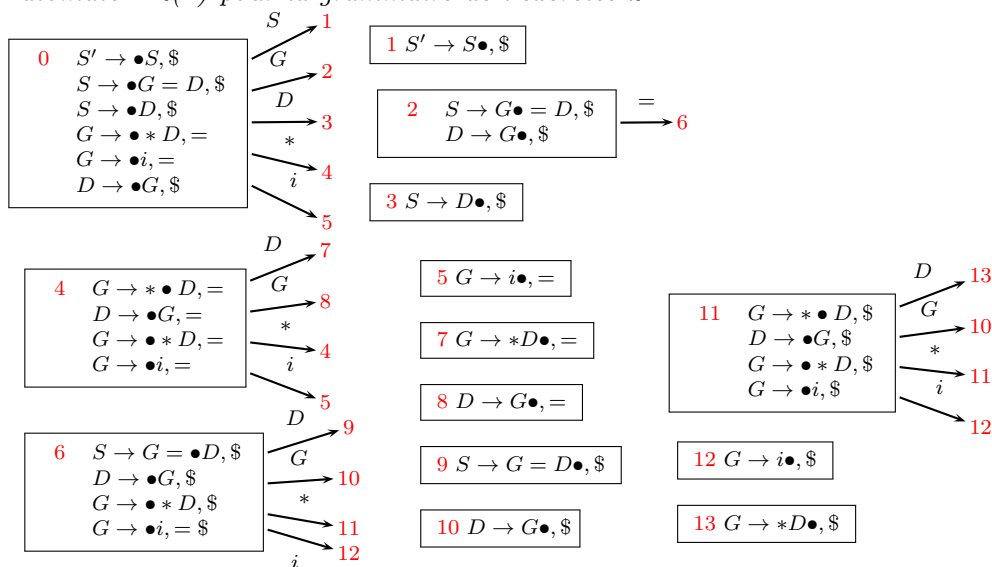
Les états Q de l'automate LR(1) seront des ensembles saturés de règles pointées étendues. Ils sont souvent appelés *collections LR(1)* (ou *LR(1)-items*). L'état initial de l'automate est $\text{Saturation}(\{S' \rightarrow \bullet S, \$\})$.

Transitions de l'automate LR(1)

Pour un ensemble saturé de règles pointées étendues E et $x \in \Sigma \cup V$, on définit l'état $\delta(E, x)$ comme

$$\text{Saturation}(\{X \rightarrow \beta x \bullet \gamma, a \mid X \rightarrow \beta \bullet x \gamma, a \in E\})$$

Automate LR(1) pour la grammaire de l'exercice 2



Remplissage des tables à partir de l'automate LR(1)

Table Successeur : pour tout $q \in Q$ et $X \in V$, si $\delta(q, X) = q'$ alors mettre q' dans la case (q, X) (identique aux méthodes LR(0) et SLR(1)).

Table Action :

- pour tout $a \in \Sigma$, $q \in Q$, si $\delta(q, a) = q'$ alors mettre “ dq' ” dans la case (q, a)
- pour tout $a \in (\Sigma \cup \{\$\})$, $q \in Q$, si q contient une règle pointée complète $X \rightarrow \beta\bullet, a$ avec $X \neq S'$, alors mettre “ $r X \rightarrow \beta$ ” dans la case (q, a)
- mettre “accepter” dans la case $(q, \$)$ où q est l'état contenant $S' \rightarrow S\bullet, \$$.
- mettre “erreur” dans toutes les cases encore vide.

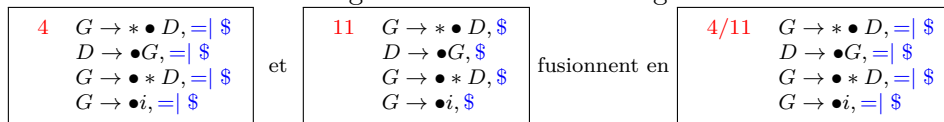
Table Action pour la grammaire de l'exercice 2

	*	i	=	\$		*	i	=	\$
0	d4	d5			7			rG → *D	
1				accepter	8			rD → G	
2			d6	rD → G	9			rS → G = D	
3				rS → D	10			rD → G	
4	d4	d5			11	d11	d12		
5			rG → i		12			rG → i	
6	d11	d12			13			rG → *D	

L'automate LR(0) comportait 10 états, l'automate LR(1) comporte 14 états. De manière générale, les automates LR(1) peuvent comporter des milliers d'états (produit cartésien) pour des langages de programmation classiques.

Automates LALR(1)

Dans l'automate LR(1) décrit précédemment, il existe des états qui ne diffèrent d'autres états que par les symboles de prévision associés aux règles pointées. Lorsque deux états de l'automate LR(1) ne diffèrent que sur les caractères de prévision, on dit qu'ils ont le même cœur LR(0). L'idée de la méthode LALR(1) est de fusionner les états LR(1) ayant le même cœur LR(0) en créant un état dont les règles sont l'union des règles des deux états fusionnés. Par exemple,

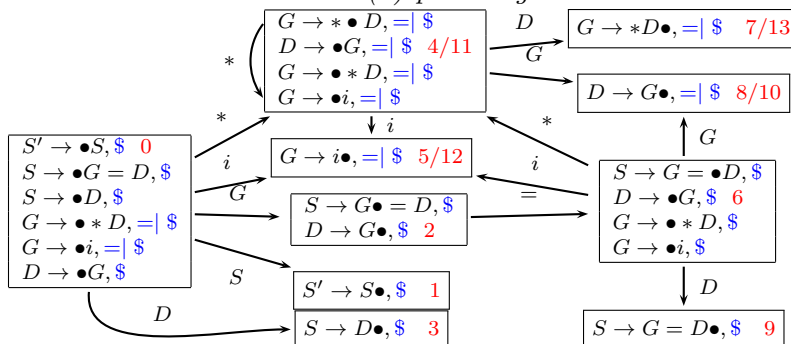


On construit l'automate LALR(1) à partir de l'automate LR(1) en fusionnant les états ayant le même cœur LR(0). Les transitions de l'automate sont simplement "mises à jour" en fonction des fusions effectuées. Remarquez que dans l'automate LR(1), les cœurs LR(0) des états sont précisément les états de l'automates LR(0). Donc l'automate LALR(1) a exactement le même nombre d'états que l'automate LR(0).

Le remplissage des tables Action et Suivant LALR(1) se fait comme le remplissage des tables LR(1) mais à partir de l'automate LALR(1). Si la table Action ne comporte qu'une action par case, alors la grammaire est LALR(1). Voici quelques remarques à noter :

- toute grammaire LALR(1) est LR(1).
- la fusion des états peut créer des conflits qui n'existaient pas dans l'automate LR(1)
 ⇒ Il existe des grammaires LR(1) qui ne sont pas LALR(1).
- les grammaires LR(1) et donc, LALR(1) ne sont pas ambiguës.

Automate LALR(1) pour la grammaire de l'exercice 2 Table Action LALR(1)



	*	i	=	\$
0	d4/11	d5/12		
1				accepter
2			d6	rD → G
3				rS → D
4/11	d4/11	d5/12		
5/12			rG → i	rG → i
6	d4/11	d5/12		
7/13			rG → *D	rG → *D
8/10			rD → G	rD → G
9				rS → G = D