



next up previous

Analyse Syntaxique ascendante

Définition d'analyse ascendante LR
 Définition de poignée et préfixe viable
 Structure d'un analyseur ascendant
 Définition et construction des Items LR(k)
 Définition et construction des tables LR(k)
 Un exemple LR(0)
 Un exemple LR(1)
 Construction des tables SLR
 Construction des tables LALR(1)

Analyse ascendante

Cette méthode cherche à construire une dérivation droite en la reconstituant à l'envers à partir de la chaîne de terminaux vers le symbole initial, ou de façon équivalente, cherche à construire un arbre de dérivation à partir des feuilles selon un parcours inverse d'un parcours en profondeur d'abord de gauche à droite.

Pour savoir simplement quand l'analyse s'arrête avec succès pour une grammaire G , on travaille toujours sur une grammaire dite *augmentée* G qui est G avec un nouveau symbole s comme symbole de départ, un nouveau symbole terminal $\$$ et une transition supplémentaire

•

Terminologie

Notation 1.1 Dans ce qui suit, on utilisera les conventions suivantes:

des lettres majuscules

indiquent des symboles non-terminaux

des lettres grecques A, \dots, Z

indiquent des *séquences* de symboles terminaux ou non-terminaux

des lettres minuscules $\alpha, \beta, \gamma, \dots$

indiquent des symboles terminaux

des lettres minuscules a, b, c, \dots

indiquent des *séquences* de symboles terminaux

Définition 1.2 Analyseurs LR Les analyseurs ascendants le plus connus sont dans la classe LR des analyseurs qui lisent le flot de tokens en entrée *de gauche à droite* (le L dans LR) pour reconstruire une dérivation *droite* (le R dans LR).

Définitions techniques

Définition 1.3 (Dérivation droite) Une dérivation droite est une dérivation qui remplace à chaque étape le symbole non terminal le

plus à droite.

On notera $\dots, u, v,$ une étape de dérivation droite entre $\alpha =$ et α , et on notera β une série (même vide) d'étapes de

dérivation droite entre $\alpha =$ et α .

Définition 1.4 (Protphrase d'une grammaire $\text{tex2html_wrap_inline}\$G\$$) Une protphrase est une séquence de symboles terminaux et non terminaux qui peut apparaître en cours d'une dérivation du symbole initial S d'une grammaire G .
On parle de protphrase *droite* (resp. *gauche*) lorsque cette séquence peut apparaître dans une dérivation droite (resp. gauche) de G .

Définitions techniques, suite

Définition 1.5 (Poignée (handle)) Dans une protphrase $\alpha =$, la séquence ϕ est une *poignée* à la position n pour la grammaire G si elle est la partie gauche d'une production γ , et que cette production doit être appliquée à $\alpha =$ en position n pour construire la protphrase précédente dans une dérivation droite à partir de S vers $\alpha =$ avec la grammaire G .

Définition 1.6 (Préfixe viable) Une séquence ϕ est un *préfixe viable* pour une grammaire G si ϕ est un préfixe de $X \rightarrow$, où

$\alpha\beta$ est une protphrase droite de G est α est une poignée dans cette protphrase.

Autrement dit, un préfixe viable est un préfixe ϕ d'une protphrase $\alpha =$, mais qui ne s'étend pas plus à droite d'une poignée α de $\alpha =$.

Un exemple

Sur la grammaire augmentée

$$\phi = \alpha\beta w$$

l'unique dérivation droite pour est la suivante:

$$id + id + id \$$$

Chaque ligne est une protphrase droite, en surligné les symboles produits, et en souligne les poignées.

Les préfixes viables sont les préfixes qui ne s'étendent pas plus loin qu'une poignée. Par

exemple, sur la protphrase ce sont $E + \underline{id} + id \$$.

Définitions techniques, fin

Définition 1.7 ($\text{tex2html_wrap_inline}\$FIRST_k()\$$) Étant donnée une grammaire G , l'ensemble $\epsilon, E, E+, E+$ contient les préfixes de longueur k des séquences de non terminaux de longueur au moins k dérivables à partir de ϕ dans G , et les séquences de non terminaux de longueur inférieur à k dérivables depuis ϕ .

Définition 1.8 ($\text{tex2html_wrap_inline}\$EFF_k()\$$ ($\text{tex2html_wrap_inline}\$-free\ FIRST_k()\$$)) Étant donnée une grammaire G , EFF_k est le sousensemble de $FIRST_k$ obtenu en considérant seulement les dérivations qui ne réduisent pas sur F un non terminal en tête de chaîne.

Exemples

Pour la grammaire:

$$\epsilon$$

on a:

$$FIRST_2(S)$$

$$\begin{aligned}
 &= \left| \begin{array}{c} \\ \\ \\ \end{array} \right| \\
 EFF_2(S) &= \{\epsilon, a, b, c,
 \end{aligned}$$

Grammaires LR(k)

On peut maintenant donner la définition formelle de la classe des grammaires $LR(k)$.

Définition 1.9 (Grammaire LR(k)) Une grammaire G est dans $LR(k)$ ($\{ca, cb\}$) si les trois conditions suivantes:

$$k \geq 0$$

$$S \Rightarrow_a^* \alpha A w \Rightarrow_a \alpha \beta w$$

$$FIRST_k(w) = FIRST_k(y)$$

$$\text{impliquent } S \Rightarrow_a^* \gamma B x \Rightarrow_a \alpha \beta y$$

Autrement dit, il suffit de regarder les premiers k symboles en entrée au moment où il faut choisir une production pour la réduction dans la dérivation droite.

Structure de l'analyseur

Les grammaires $LR(k)$ sont celles dont le langage est reconnu par un analyseur déterministe $LR(k)$. Cet analyseur utilise une pile et le flot d'entrée, qui décrivent une *configuration* de l'analyseur, notée

$$\alpha = \gamma, A = B, x = y$$

où les x sont des symboles, terminaux ou non terminaux, stockés sur la pile, alors que les a sont seulement des symboles terminaux, et correspondent aux terminaux non encore lus sur le flot d'entrée.

L'analyseur travaille en effectuant quatre actions possibles:

shift (décalage)

on transfère le terminal a_i du flot d'entrée vers la pile

reduce (réduction)

on reconnaît sur le sommet de la pile une partie droite d'une production

, on l'enlève et on la

remplace par sa partie gauche Y

erreur

l'analyseur s'arrête et signale une erreur

accept

l'analyseur s'arrête et signale que la phrase a été reconnue

Pour choisir les actions, on utilise une table d'analyse que l'on verra plus avant.

Un exemple

Sur la grammaire augmentée

$$\phi = \alpha \beta w$$

Une possible séquence de reconnaissance pour

pour un analyseur

ascendant serait:

Remarques Importantes

la concaténation de la partie gauche et droite d'une configuration d'un analyseur ascendant pour une grammaire G est toujours une protophrase *droite* de G (si l'analyse se termine avec succès).

un préfixe viable peut toujours se compléter en une protophrase droite. En d'autre termes, il n'y a pas d'erreur au cours de l'analyse tant que l'on a sur la pile un préfixe viable.

Un autre exemple, avec look-ahead

(,	$id + id + id$	\$)	<i>shift</i>
(<u>id</u>	,	$+id + id$	\$)	<i>reduce</i>
(<u>T</u>	,	$+id + id$	\$)	<i>reduce</i>
(<u>E</u>	,	$+id + id$	\$)	<i>shift</i>

ascendant serait:

$$\begin{array}{lcl} S & \rightarrow & E \$ \\ E & \rightarrow & T + E \mid T \\ T & id & \end{array}$$

Analyseurs LR

Un analyseur LR est composé de

qui décrit un automate à états finis augmenté avec des actions à effectuer éventuellement sur la pile (shift, reduce, accept, error)

L'exécution de l'automate est censée décaler sur la pile des symboles jusqu'à atteindre une préfixe viable maximale (i.e. pas extensible à droite, i.e. contenant une poignée ϕ en fond à droite, i.e. en sommet de pile), puis réduire la poignée en la remplaçant avec la partie droite x de la production concernée.

Fonctionnement d'un analyseur LR

Sur un état d'analyseur $x \rightarrow \gamma$ le fonctionnement de l'analyseur LR est le suivant:

shift (noté s) déplacer le symbole d'entrée x sur la pile.

(noté rn) sur le sommet de la pile il y a la partie gauche de la règle numéro n , disons γ ; dépiler ϕ et empiler X

accept

(noté a) arrêter avec succès

error

(noté par une case vide) arrêter sur erreur

recommencer avec le nouvel état d'analyseur

Exemple d'exécution avec une table d'analyse LALR(1)

```
<!-- MATH: \begin{displaymath} \begin{array}{c@{\quad}c@{\quad}\rightarrow\quad}@{\quad}c@{\quad}c@{\quad}\rightarrow\quad}\{ 0 \& S \& E; \& \$ \& 2 \& E \& T \setminus 1 \& E \& T + E \& 3 \& T \& id \end{array} \\ \end{displaymath} -->
```

$$(\alpha, \quad xw)$$

0	S	\rightarrow	$E \$$	2	E	\rightarrow	T
1	E	\rightarrow	$T + E$	3	T	\rightarrow	id

Ici on a marqué en bas les états de l'automate après lecture de chaque symbole sur la pile.

Analyseur avec états sur la pile

Si on garde les états sur la pile, en modifiant la notion de configuration pour que chaque symbole soit suivi par un état, on peut éviter de relire toute la pile à chaque fois: sur une configuration d'analyseur

	<i>Action</i>	<i>Transition</i>
	$id + \$$	$id + \$ E \uparrow$

le fonctionnement de l'analyseur LR est le suivant:

exécuter l'action décrite dans la table d'analyse associée au symbole terminal x en entrée pour l'état s_k

shift k

déplacer le symbole d'entrée x sur la pile, et empiler l'état numéro k

reduce n

sur le sommet de la pile il y a la partie gauche de la règle numéro n , disons γ ; dépiler ϕ et tous les états associés,

en découvrant l'état s' ; empiler X et l'état s' contenu dans la table à la ligne s' , colonne X

accept

arrêter avec succès

error

arrêter sur erreur

recommencer avec le nouvel état d'analyseur

Comment produire une table d'analyse?

Il faut savoir reconnaître les préfixes viables, et savoir déterminer quelles productions utiliser pour les réductions, éventuellement en utilisant k tokens en entrée pour aider dans la décision.

Pour reconnaître les préfixes viables, on définit d'abord

Définition 1.10 (ITEM LR(k)) Un *ITEM LR(k)* pour une grammaire G est une production

de G plus une position j dans ϕ et une séquence w de longueur $(s_1 X_1)$. Cela est noté, si $\leq k$ avec j la

longueur de $\alpha =$

$$\gamma = \alpha\beta$$

sauf dans le cas LR(0) pour lequel on écrit simplement

|

L'intuition de $X \rightarrow \alpha \cdot \beta$ est que l'on a déjà vu en entrée le préfixe α d'une protophrase et que l'on attende sur l'entrée une séquence dérivable à partir de β .

Reconnaître les préfixes viables: la fermeture

Si on a βw , i.e. on a déjà vu en entrée le préfixe α et on attende une séquence dérivable à partir de β , on est aussi en condition d'attendre une séquence dérivable depuis α , suivie d'une séquence dérivable depuis β . C'est cela que capture la notion suivante de fermeture

Définition 1.11 (Fermeture (Closure) LR(k))

Fermeture(I) =
répéter tant que I grandit
pour tout item βw dans I

pour toute production βz pour tout $X \rightarrow \gamma w \in FIRST_k(\beta z)$ retourner I

Reconnaître les préfixes viables: GOTO

Définition 1.12 (GOTO)

Supposons d'avoir βw , pour un symbole *terminal ou non terminal* X : on a donc déjà vu en entrée le préfixe α et on attende une séquence dérivable à partir de β . Si maintenant l'on reconnaît X , alors on a vu αX et on attende une séquence dérivable à partir de β .

C'est cela que capture la notion suivante de GOTO

Goto(I, X) =
 αX pour tout item βw dans I
 $J \leftarrow \emptyset$ retourner Fermeture(J)

L'automate qui reconnaît les préfixes viables

Soit G un grammaire augmentée, et soit \mathcal{I} la collection d'ensembles d'ITEMS LR(k) atteignables depuis la fermeture de l'item $\{s_0; s_1; \dots; s_k\}$ par la fonction GOTO.

On peut alors construire l'automate à état fini suivant:

états \mathcal{I} avec s_0 état initial et comme états finaux ceux qui contiennent au moins un ITEM LR(k) avec le point au fond à droite (i.e. de la forme $s_0 = (S' \rightarrow \cdot)$)

transitions on a une transition de l'état s_i vers l'état s_j sur le symbole X si $GOTO(s_i, X) = s_j$

La construction de la table LR(k)

Soit G un grammaire augmentée, pour laquelle on a construit l'automate.

La table d'analyse a une ligne par état et un colonne par séquence de symboles terminaux de longueur $(s_1 X_1)$ (le look-ahead) et une colonne par symbole terminal et non-terminal, que l'on remplit de la façon suivante:

Pour tout état $(X \rightarrow \gamma,$

on met $s_i \in \mathcal{I}$ dans la case s_i, u si *reduce* n et $(A \rightarrow \beta \cdot$ est la production numéro $A \rightarrow \beta$
 on met *accept* dans la case $n \geq 1$ si $s_i, \$$
 on met *shift* dans la case s_i, u si $(S' \rightarrow \beta \cdot, \$) \in s_i$ et $(A \rightarrow \beta_1 \cdot \beta_2, v) \in$
 on laisse vide (i.e. on signale erreur) autrement

Theorem 1.13 (fondamental de l'analyse ascendante) Si une grammaire G est LR(k), alors l'automate construit reconnaît les préfixes viables de G et tout état contenant un ITEM LR(k) de la forme $s_0 = (S' \rightarrow \cdot$ ne contient pas un ITEM $u \in EFF_k(\beta_2 v)$

avec $(X' \rightarrow \gamma_1 \cdot \gamma_2, u)$. (en d'autre terme, on n'aura pas dans la table des entrées multiples décaler et réduire ou entre deux réductions).

Comment l'analyseur peut-il choisir l'action à effectuer?

Un analyseur LR dispose de plus d'information qu'un analyseur LL pour déterminer la prochaine action.

Imaginons d'avoir en entrée une chaîne uvw , et d'avoir déjà lu u .
 Pour déterminer la production à appliquer

un analyseur LL(k) connaît u et $FIRST_k(vw)$

un analyseur LR(k) connaît uv (en effet, il connaît un préfixe viable ϕ obtenu à partir de uv) et $FIRST_k(w)$

Un exemple LR(0)

La grammaire $w \in \mathcal{I}$ suivante est LR(0)

$$G \quad (1)$$

Voici la construction complète de la table LR(0) de G
 États et transitions (2,4,6,7,9 sont terminaux)

1. $G \quad 1$
 $\{(S' \rightarrow \cdot S \$), (S \rightarrow \cdot (L)), (S \rightarrow \cdot x)\}$
2. $\text{goto}(1, \$$
3. $\{(S \rightarrow x \cdot)\}$
 $\{(S \rightarrow \cdot (L)), (L \rightarrow \cdot S), (L \rightarrow \cdot L, S), (L \rightarrow \cdot S)(S \rightarrow \cdot (L)), (S \rightarrow \cdot x)\}$
4. $\text{goto}(3, () =$
5. $\{(S' \rightarrow S \cdot \$)\}$
 $\{(S \rightarrow (L \cdot)), (L \rightarrow L \cdot, S)\}$
6. $\text{goto}(5,)) :$
7. $\{(S \rightarrow (L) \cdot)\}$
8. $\{(L \rightarrow S \cdot)\}$
 $\{(L \rightarrow L, \cdot S), (S \rightarrow \cdot (L)), (S \rightarrow \cdot x)\}$
- 9.

goto(8, x) =

La table d'analyse LR(0)

Pour remplir la table LR(0), on écrit la table de transition de l'automate et on introduit les actions de décalage (s pour $shift$) comme décrit plus en haut.

Pour les réductions (r_k pour $reduce$ avec la production k), n'ayant pas de look-ahead dans les états, on met r_k dans toute la ligne action de l'état j si l'état j contient un ITEM LR(0) $\{(L \rightarrow L, \text{ et que } \gamma \text{ est la production numéro } j.$

$$X \rightarrow \gamma.$$

La table d'analyse LR(0), versions compacte

Remarque

Les générateurs d'analyseurs, comme Yacc, fusionnent les colonnes des actions et des transitions pour les terminaux: plutôt que d'avoir une case case $(1, x)$ qui contient s (hift) pour les actions, et une case $(1, x)$ qui contient 2 pour les transitions, on préfère avoir une seule case $(1, x)$ qui contient $s2$, pour l'action est un shift et la transition est vers 2. Dans ce cas, on écrit souvent $s2$ dans les colonnes transitions restantes (celles des non-terminaux). C'est une abréviation pour $goto$ k , plus lisible que juste k .

	Action					Transition				
	()	x	,	$\$$	()	x	,	$\$$
1	s		s			3		2		4
2	$r2$	$r2$	$r2$	$r2$	$r2$					
3	s		s			3		2		7 5
4					a					
5		s		s			6		8	
6	$r1$	$r1$	$r1$	$r1$	$r1$					
7	$r3$	$r3$	$r3$	$r3$	$r3$					
8	s		s			3		2		9
9	$r4$	$r4$	$r4$	$r4$	$r4$					

Un exemple LR(1) non LR(0)

La grammaire

$$(\alpha \rightarrow xw) \quad (2)$$

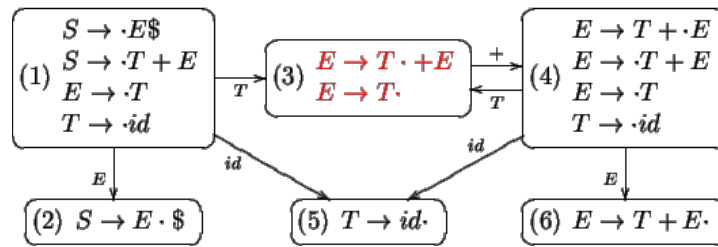
est une grammaire $LR(1)$ mais pas $LR(0)$.

En effet, dans l'automate on a un problème pour l'état G_2 .

$$\{(E \rightarrow T \cdot +E), (E \rightarrow T \cdot)\}$$

Un exemple LR(1) non LR(0) (suite)

Donc dans la table d'analyse $LR(0)$ on trouve un conflit $shift/reduce$ dans la case $3,+$



Les états LR(1) de

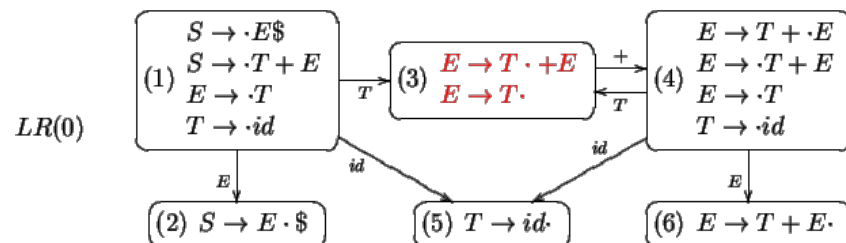
Regardons alors la construction LR(1), qui garde trace des look-aheads dans les états...

États et transitions (2, 5 et 6 sont terminaux)

G_2

Comparons les automates LR(0) et LR(1) pour

G_2



La table d'analyse LR(1) de

On garde trace des look-ahead pour introduire dans la table les actions *reduce*, donc il y en a moins et le conflit disparaît!

G_2

La table d'analyse LR(1) de, version compacte

G_2

Comparons les tables LR(0) et LR(1) pour

G_2

Trouver sa place parmi les LR(k)

Comme nous venons de voir, la classe d'analyseurs LR(0) est trop faible pour traiter les langages de programmations: même le simple langage des expressions pose problème.

Les classes LR(2), LR(3), ... ont par contre une table d'analyse trop grosse en pratique en raison du nombre de colonnes pour le "look-ahead": un analyseur moderne utilise plusieurs dizaines de tokens, et une colonne pour chaque séquence de token de longueur inférieur ou égale à k , pour k est déraisonnable.

Exercice: combien de séquences de longueur inférieur ou égale à k y-at-il si on se donne n tokens différents?

Trouver sa place entre LR(0) et LR(1)

Heureusement, la classe LR(1) est largement suffisante pour les langages modernes, et la table n'a qu'une colonne par token. Mais là, c'est le nombre d'états qui grandit trop, en raisons de la présence de look-ahead dans les états qui départage des états qui sont très peu différents.

C'est pour cela que dans la pratique on utilise deux types d'analyseurs dont la puissance est comprise entre celle de LR(0) et celle de LR(1): SLR et LALR(1)

Analyseurs SLR

SLR

(Simple LR) est un analyseur dont l'automate est celui de LR(0), donc la partie transition est la même que LR(0), et les actions de décalage aussi, mais la table d'analyse est construite de une façon plus fine: on pallie à l'absence de look-ahead dans les état avec l'information contenue dans les ensembles FOLLOW construits à partir de la grammaire. La règle de placement des réductions devient alors:

si l'état j contient un ITEM LR(0) $\{L \rightarrow L_j$, et que γ est la production numéro $k \geq 2$, on met rk dans

toutes les cases (j, t) telles que $t \in FOLLOW(X)$.

SLR pour la grammaire G_2

L'automate SLR étant le même que celui LR(0), on ne le montrera pas à nouveau, mais maintenant la table SLR contiendra des entrées *reduce* seulement sur certains nonterminaux, pas tous! En particulier, on pourra éviter le conflit *shift/reduce* dans l'état G_2 .

$t \in FOLLOW(X)$

Dans ce cas précis, SLR fait aussi bien que LR(1), avec bien moins d'effort.

Analyseurs LALR(1)

LALR(1)

(Look-Ahead LR(1)) est une classe d'analyseurs dont l'automate est obtenu de l'automate LR(1) en fusionnant les états qui diffèrent seulement par leur look-ahead.

On dit aussi que l'on fusionne les états ayant le même *coeur*, le coeur d'un état étant l'ensemble des parties gauches des ITEMS LR(1) qu'il contient, i.e. sans le look-ahead, i.e. des ITEMS LR(0). Donc un analyseur LALR(1) a autant d'états qu'un LR(0) ou SLR.

Les analyseurs LALR(1) sont les plus utilisés parce que, même s'ils ont moins d'états qu'un analyseur LR(1), il est très rare qu'on retrouve un conflit dans la table LALR(1) quand il n'y en a pas dans la table LR(1).

En particulier, on peut prouver que si un analyseur LR(1) n'a pas de conflits *shift/reduce*, l'analyseur LALR(1) n'en a pas non plus. Par contre, on peut introduire des conflits *reduce/reduce*.

LALR(1) pour G_2

Dans le cas précis de cette grammaire, l'automate LR(1) pour G_2 n'ayant pas d'états différents avec le même coeur, la table d'analyse LALR(1) de G_2 est la même que celle LR(1).

Mais la grammaire suivante, qui capture un sous-ensemble des expressions du langage C, est un exemple de grammaire LALR(1) qui n'est pas SLR et pour laquelle l'automate LALR(1) est plus petit que l'automate LR(1).

G_2

(3)

LR préfère l'associativité à gauche

Contrairement à ce qui se passe dans le cas des analyseurs LL, dans l'analyse ascendante on a plutôt intérêt à utiliser des grammaires récursives à gauche.

Considérons les analyseurs LR pour la grammaire récursive à droite

0	S'	\rightarrow	$S \$$	3	E	\rightarrow	V
1	S	\rightarrow	$V = E$	4	V	\rightarrow	id
2	E	\rightarrow	$E + E$	5	E	\rightarrow	$E * E$

vus en cours: pour reconnaître $id + id + \dots + id \$$, ils empilent

toute la suite de symboles (en réduisant id sur τ à chaque coup) avant de faire la première réduction non triviale.

Par contre, la grammaire récursive à gauche

$id + id + \dots + id \$$

mantiendra la dimension de la pile à un minimum.

Utilisation de grammaires ambiguës

Une grammaire ambiguë n'est jamais LR(k), quelque soit k .
 Pourtant, on a intérêt à essayer d'utiliser une grammaire ambiguë, quitte à trafiquer l'automate $LR(k)$, si on peut.

efficacité

dans une grammaire obtenue par désambiguation, l'analyseur passe beaucoup de temps à réduire des productions triviales

(comme | dans l'exemple précédent), dont le seul but était d'explicitier *dans la grammaire* les priorités entre opérateurs et leur associativité droite ou gauche.

praticité

si on peut décrire de façon concise ces priorités entre opérateurs et leur associativité droite ou gauche, sans toucher à la grammaire, on obtient une description plus modulaire du langage qui nous intéresse.

Exemple

Considérons la grammaire (ambiguë) suivante:

```
<!-- MATH: \begin{displaymath} \begin{array}{c} \rightarrow \\ S \ \& \ E; \ \$ \ E \ \& \ E * E \ \& \ E + E \ \& \ id \end{array} \end{displaymath} -->
```

$$E \rightarrow T$$

et sa table d'analyse SLR.

Voyons comment les nombreux conflits apparents peuvent s'expliquer en terme d'associativité et précédence d'opérateurs, que l'on peut résoudre *en travaillant directement sur les entrées de la table...*

(fait au tableau, pas dans les notes... si un ame gentille veut tout taper...)

About this document ...

This document was generated using the [LaTeX2HTML](#) translator Version 98.1p1 release (March 2nd, 1998)

Copyright © 1993, 1994, 1995, 1996, 1997, [Nikos Drakos](#), Computer Based Learning Unit, University of Leeds.

The command line arguments were:
`latex2html -split 0 Slides01.`

The translation was initiated by Roberto Di Cosmo on 1999-11-24

next up previous

Roberto Di Cosmo
 1999-11-24