

Notes on LR parsing

Handout written by Maggie Johnson and revised by me.

The problem with SLR

In SLR(1), we do not use all the information that we have at our disposal. When we have a completed configuration (i.e. dot at the end) such as $X \rightarrow \cdot$, we know this corresponds to a situation in which we have \underline{u} as a handle on top of the stack which we then can reduce, i.e., replacing \underline{u} by X . We allow such a reduction whenever the next symbol is in $\text{Follow}(X)$. However, it may be that we should not reduce for every symbol in $\text{Follow}(X)$, because the symbols below \underline{u} on the stack preclude \underline{u} being a handle for reduction in this case. In other words, SLR(1) states only tell us about the sequence on top of the stack, not what is below it on the stack. We may need to divide an SLR(1) state into separate states to differentiate the possible means by which that sequence has appeared on the stack. By carrying more information in the state, it will allow us to rule out these invalid reductions.

At the end of the SLR handout, we give the following example:

$S' \rightarrow S$
 $S \rightarrow L = R$
 $S \rightarrow R$
 $L \rightarrow *R$
 $L \rightarrow id$
 $R \rightarrow L$

$I_0: S' \rightarrow S$
 $S \rightarrow L = R$
 $S \rightarrow R$
 $L \rightarrow *R$
 $L \rightarrow id$
 $R \rightarrow L$

$I_1: S' \rightarrow \cdot S$

$I_2: S \rightarrow L \cdot = R$
 $R \rightarrow \cdot$

$I_3: S \rightarrow R \cdot$

$I_4: L \rightarrow \cdot *R$
 $R \rightarrow \cdot L$
 $L \rightarrow \cdot id$

$I_5: L \rightarrow id \cdot$

$I_6: S \rightarrow L = \cdot R$
 $R \rightarrow \cdot$
 $L \rightarrow \cdot *R$
 $L \rightarrow \cdot id$

$I_7: L \rightarrow *R \cdot$

$I_8: R \rightarrow L \cdot$

$I_9: S \rightarrow L = R \cdot$

Consider parsing the expression $id = id$. After working our way to configuring set I_2 having reduced the first id to L , we have a choice upon seeing $=$ coming up in the input. The first configuration in the set wants to set $\text{Action}[2,=]$ be shift 6, which corresponds to moving on to find the rest of the assignment. However, $=$ is also in $\text{Follow}(R)$ because $S \Rightarrow L=R \Rightarrow *R = R$. Thus, the second configuration wants to reduce in that slot $R \rightarrow \cdot$. This is a shift-reduce conflict but not

because of any problem with the grammar. A SLR parser does not remember enough left context to decide what should happen when it encounters a $=$ in the input having seen a string reducible to L . Although the sequence on top of the stack could be reduced to R , but we don't want to choose this reduction because there is no possible right sentential form that begins $R = \dots$ (there is one beginning $*R = \dots$ which is not the same). Thus, the correct choice is to shift.

It's not further lookahead that the SLR tables are missing, we don't need to see additional symbols coming up in the input, we have already seen the information that allows us to determine the correct choice. What we need is to retain a little more of the left context that brought us here. In this example grammar, the only time we should consider reducing by production $R \rightarrow \epsilon$ is during a derivation that has already seen a $*$ or an $=$. Just using the entire follow set is not discriminating enough as the guide for when to reduce. The follow set contains symbols that can follow R in any valid sentence but it does not precisely indicate which symbols follow R within a particular derivation. So we will augment our states to include information about exactly what portion of the follow set is appropriate given the path we have taken to that state.

We can be in state 2 for one of two reasons, we are trying to build a sentence from $S \rightarrow \epsilon = R$ or from $S \rightarrow R \rightarrow \epsilon$. If the upcoming symbol is $=$, then that rules out the second choice and we must be building the first, which tells us to shift. The reduction should only be applied if the next input symbol is $\$$. Even though $=$ is $\text{Follow}(R)$ because of the other contexts that an R can appear, in this particular situation, it is not appropriate because when deriving a sentence $S \rightarrow R \rightarrow \epsilon$, $=$ cannot follow R .

Constructing LR(1) parsing tables

LR or *canonical LR* parsing incorporates the required extra information into the state by redefining configurations to include a terminal symbol as an added component. LR configurations have the general form:

$$A \rightarrow \epsilon_1 \dots X_i \cdot X_{i+1} \dots X_j, a$$

This means we have states corresponding to $X_1 \dots X_i$ on the stack and we are looking to put states corresponding to $X_{i+1} \dots X_j$ on the stack and then reduce, but only if the token following X_j is the terminal a . a is called the *lookahead* of the configuration. The lookahead only comes into play with LR(1) configurations with a dot at the right end:

$$A \rightarrow \epsilon_1 \dots X_j \cdot, a$$

This means we have states corresponding to $X_1 \dots X_j$ on the stack but we may only reduce when the next symbol is a . The symbol a is either a terminal or $\$$ (end of input marker). With SLR(1) parsing, we would reduce if the next token was any of those in $\text{Follow}(A)$. With LR(1) parsing, we reduce only if the next token is exactly a . We may have more than one symbol in the lookahead for the configuration, as a convenience, we list those symbols separated by a forward slash. Thus, the configuration $A \rightarrow \epsilon \cdot, a/b/c$ says that it is valid to reduce ϵ to A only if the next token is equal to a , b , or c . The configuration lookahead will always be a subset of $\text{Follow}(A)$, but it could be a proper subset.

Recall the definition of a *viable prefix* from the SLR handout. Viable prefixes are those prefixes of right sentential forms that can appear on the stack of a shift-reduce parser. Formally we say that a configuration $[A \rightarrow \epsilon \cdot \underline{v}, a]$ is valid for a viable prefix \underline{v} if there is a rightmost derivation $S \Rightarrow^* A \underline{w} \Rightarrow^* \underline{uvw}$ where $\underline{v} = \underline{u}$ and either a is the first symbol of \underline{w} or \underline{w} is ϵ and a is $\$$. For example,

$Z \leftarrow Z \mid y$

There is a rightmost derivation $S \Rightarrow^* xxZxy \Rightarrow xxxZxy$. We see that configuration $[Z \leftarrow \bullet Z, x]$ is valid for viable prefix $= xxx$ by letting $= xx$, $A = Z$, $\underline{w} = xy$, $\underline{u} = x$ and $\underline{v} = Z$. Another example is from the rightmost derivation $S \Rightarrow^* ZxZ \Rightarrow ZxxZ$, making $[Z \leftarrow \bullet Z, \$]$ valid for viable prefix Zxx .

Often we have a number of LR(1) configurations that differ only in their lookahead components. The addition of a lookahead component to LR(1) configurations allows us to make parsing decisions beyond the capability of SLR(1) parsers. There is, however, a big price to be paid. There will be more distinct configurations and thus many more possible configuring sets. This increases the size of the goto and action tables considerably. In the past when memory was smaller, it was difficult to find storage-efficient ways of representing these tables, but now this is not as much of an issue. Still, it's a big job building LR tables for any substantial grammar by hand.

The method for constructing the configuring sets of LR(1) configurations is essentially the same as for SLR, but there are some changes in the closure and successor operations because we must respect the configuration lookahead. To compute the closure of a LR configuring set I:

Repeat the following until no more configurations can be added to state I:

- For each configuration $[A \leftarrow \bullet B\underline{v}, a]$ in I, for each production $B \leftarrow \underline{w}$ in G' , and for each terminal b in $\text{First}(\underline{v}a)$ such that $[B \leftarrow \underline{w}, b]$ is not in I: add $[B \leftarrow \underline{w}, b]$ to I.

What does this mean? We have a configuration with the dot before the non-terminal B. In SLR, we computed the closure by adding all B productions with no indication of what was expected to follow them. In LR, we are a little more precise—we add each B production but insist that each have a lookahead of $\underline{v}a$. The lookahead will be $\text{First}(\underline{v}a)$ since this is what follows B in this production. Remember that we can compute first sets not just for a single non-terminal, but also a sequence of terminal and non-terminals. $\text{First}(\underline{v}a)$ includes the first set of the first symbol of \underline{v} and then if that symbol is nullable, we include first set of the following symbol, and so on. If the entire sequence \underline{v} is nullable, we add the lookahead already required by this configuration.

The successor function for the configuring set I and symbol X is computed as this:

Let J be the configuring set $[A \leftarrow \underline{X}\bullet\underline{v}, a]$ such that $[A \leftarrow \underline{X}\bullet\underline{v}, a]$ is in I. $\text{successor}(I, X)$ is the closure of configuring set J.

We take each production in a configuring set, move the dot over a symbol and close on the resulting production. This is pretty much the same successor function as defined for LR, but we have to propagate the lookahead when computing the transitions.

We construct the complete family of all configuring sets F just as we did for SLR. F is initialized to the set with the closure of $[S' \leftarrow \$, \$]$. For each configuring set I and each grammar symbol X such that $\text{successor}(I, X)$ is not empty and not in F, add $\text{successor}(I, X)$ to F until no other configuring set can be added to F.

Let's consider an example. The augmented grammar below that recognizes the regular language a^*ba^*b (this example from pp. 231-236 Aho/Sethi/Ullman).

- 0) $S' \leftarrow \$$
- 1) $S \leftarrow \underline{X}$
- 2) $X \leftarrow \underline{X}$
- 3) $X \leftarrow \underline{\$}$

Here is the family of LR configuration sets:

$I_0:$	$S' \leftarrow S, \$$	$I_4:$	$X \leftarrow \bullet, a/b$
	$S \leftarrow XX, \$$	$I_5:$	$S \leftarrow X\bullet, \$$
	$X \leftarrow aX, a/b$	$I_6:$	$X \leftarrow \bullet X, \$$
	$X \leftarrow b, a/b$		$X \leftarrow aX, \$$
$I_1:$	$S' \leftarrow \bullet, \$$		$X \leftarrow b, \$$
$I_2:$	$S \leftarrow \bullet X, \$$	$I_7:$	$X \leftarrow \bullet, \$$
	$X \leftarrow aX, \$$	$I_8:$	$X \leftarrow X\bullet, a/b$
	$X \leftarrow b, \$$	$I_9:$	$X \leftarrow X\bullet, \$$
$I_3:$	$X \leftarrow \bullet X, a/b$		
	$X \leftarrow aX, a/b$		
	$X \leftarrow b, a/b$		

The above grammar would only have seven SLR states, but has ten in canonical LR. We end up with a lot more states because we have split states that have different lookaheads. For example, states 3 and 6 are the same except for lookahead, state 3 corresponds to the context where we are in the middle of parsing the first X, state 6 is the second X. Similarly, states 4 and 7 are completing the first and second X respectively. In SLR, those states are not distinguished, and if we were attempting to parse a single b by itself, we would allow that to be reduced to X, even though this will not lead to a valid sentence. The SLR parser will eventually notice the syntax error, but the LR parser figures it out a bit sooner.

To fill in the entries in the action and goto tables, we use the pretty much the same algorithm as we did for SLR. It is a bit easier to assign the reduce actions for an LR table because the lookahead sets tell us exactly where to put each reductions.

1. Construct $F = \{I_0, I_1, \dots, I_n\}$, the collection of configuring sets for the augmented grammar G' (augmented by adding the special production $S' \leftarrow \bullet$).
2. State i is determined from I_i . The parsing actions for the state are determined as follows:
 - a) If $[A \leftarrow \bullet a, b]$ is in I_i and $\text{succ}(I_i, a) = I_j$, then set $\text{Action}[i, a]$ to shift j (a must be a terminal).
 - b) If $[A \leftarrow \bullet, a]$ is in I_i then set $\text{Action}[i, a]$ to reduce $A \leftarrow \bullet$ (note A may not be S').
 - c) If $[S' \leftarrow \bullet, \$]$ is in I_i then set $\text{Action}[i, \$]$ to accept.
3. The goto transitions for state i are constructed for all non-terminals A using the rule: If $\text{succ}(I_i, A) = I_j$, then $\text{Goto}[i, A] = j$.
4. All entries not defined by rules 2 and 3 are errors.
5. The initial state is the one constructed from the configuring set containing $S' \leftarrow \bullet S$.

Following the algorithm using the configuring sets given above, we construct this canonical LR parse table:

State on top of stack	Action			Goto	
	a	b	\$	S	X
0	s3	s4		1	2
1			acc		
2	s6	s7			5
3	s3	s4			8
4	r3	r3			
5			r1		
6	s6	s7			9
7			r3		
8	r2	r2			
9			r2		

Let's parse an example string baab. It is a valid in this language as shown by this leftmost derivation:

S \leftarrow XX
 bX
 baX
 baaX
 baab

STACK	REMAINING	PARSER
STACK	INPUT	ACTION
S ₀	baab\$	Shift S ₄
S ₀ S ₄	aab\$	Reduce 3) X \rightarrow bX, goto S ₂
S ₀ S ₂	aab\$	Shift S ₆
S ₀ S ₂ S ₆	ab\$	Shift S ₆
S ₀ S ₂ S ₆ S ₆	b\$	Shift S ₇
S ₀ S ₂ S ₆ S ₆ S ₇	\$	Reduce 3) X \rightarrow bX, goto S ₉
S ₀ S ₂ S ₆ S ₆ S ₉	\$	Reduce 2) X \rightarrow baX, goto S ₉
S ₀ S ₂ S ₆ S ₉	\$	Reduce 2) X \rightarrow baX, goto S ₅
S ₀ S ₂ S ₅	\$	Reduce 1) S \rightarrow baX, goto S ₁
S ₀ S ₁	\$	Accept

Now, let's consider what the states mean. S₄ is where X \rightarrow bX is completed; S₂ and S₆ is where we are in the middle of processing the 2 a's; S₇ is where we process the final b; S₉ is where we complete the X \rightarrow baX production; S₅ is where we complete S \rightarrow baX; and S₁ is where we accept.

LR(1) grammars

Every SLR(1) grammar is a canonical LR(1) grammar, but the canonical LR parser may have more states than the SLR parser for same grammar. An LR(1) grammar may not be SLR(1), the grammar on page 1 of this handout is such an example. Those grammars require splitting states to resolve shift-reduce conflicts that would otherwise result if lookaheads were not carried with the states.

A grammar is LR(1) if the following two conditions are satisfied for each configuring set:

1. For any item in the set $[A \rightarrow \alpha \bullet x \gamma, a]$ with x a terminal, there is no item in the set of the form $[B \rightarrow \beta \bullet, x]$. In the action table, this translates no shift-reduce conflict for any state. The successor function for x either shifts to a new state or reduces, but not both.
2. The lookaheads for all complete items within the set must be disjoint, e.g. set cannot have both $[A \rightarrow \alpha \bullet, a]$ and $[B \rightarrow \beta \bullet, a]$. This translates to no reduce-reduce conflict on any state. If more than one non-terminal could be reduced from this set, it must be possible to uniquely determine which is appropriate from the next input token.

As long as there is a unique shift or reduce action on each input symbol from each state, we can parse using an LR(1) algorithm. The above state conditions are similar to what is required for SLR(1), but rather than the looser constraint about disjoint follow sets and so on, canonical LR computes a more precise notion of the appropriate lookahead within a particular context and thus is able to resolve conflicts that SLR would encounter.

Bibliography

- A. Aho, R. Sethi, J. Ullman, Compilers: Principles, Techniques, and Tools. Reading, MA: Addison-Wesley, 1986.
- J.P. Bennett, Introduction to Compiling Techniques. Berkshire, England: McGraw-Hill, 1990.
- K. Loudon, Compiler Construction. Boston, MA: PWS, 1997.
- A. Pyster, Compiler Design and Construction. New York, NY: Van Nostrand Reinhold, 1988.