

Fixing problems with grammars

Informatics 2A: Lecture 12

Alex Simpson

School of Informatics
University of Edinburgh
als@inf.ed.ac.uk

12 October, 2011

LL(1) grammars: summary

Given a context-free grammar, the problem of **parsing** a string can be seen as that of constructing a **leftmost derivation**, e.g.

$$\begin{aligned} \text{Exp} &\Rightarrow \text{Exp} + \text{Exp} \\ &\Rightarrow \text{Num} + \text{Exp} \\ &\Rightarrow 1 + \text{Exp} \\ &\Rightarrow 1 + \text{Num} \\ &\Rightarrow 1 + 2 \end{aligned}$$

At each stage, we expand the **leftmost nonterminal**. In general, it (seemingly) requires **magical powers** to know which rule to apply.

An **LL(1) grammar** is one in which the correct rule can always be determined from just the **nonterminal** to be expanded and the current **input symbol** (or end-of-input marker).

This leads to the idea of a **parse table**: a two-dimensional array (indexed by nonterminals and input symbols) in which the appropriate production can be looked up at each stage.

Possible problems with grammars

LL(1) grammars allow for very **efficient parsing** (time linear in length of input string). Unfortunately, many “natural” grammars are not LL(1), for various reasons, e.g.

- ① They may be **ambiguous** (bad for computer languages)
- ② They may have rules with **shared prefixes**: e.g. how would we choose between the following productions?

$$\text{Stmt} \rightarrow \text{do Stmt while Cond}$$

$$\text{Stmt} \rightarrow \text{do Stmt until Cond}$$

- ③ They may be **left-recursive rules**, where the LHS nonterminal appears at the start of the RHS: $\text{Exp} \rightarrow \text{Exp} + \text{Exp}$

Sometimes such problems can be fixed: we can replace our grammar by an equivalent LL(1) one. We'll look at some ways of doing this.

Problem 1: Ambiguity

We've seen many examples of **ambiguous** grammars. Some kinds of ambiguity are 'needless' and can be easily avoided. E.g. can replace

$$\text{List} \rightarrow \epsilon \mid \text{Item} \mid \text{List List}$$

by

$$\text{List} \rightarrow \epsilon \mid \text{Item List}$$

A similar trick works generally for any other kind of 'lists'.
E.g. can replace

$$\text{List1} \rightarrow \text{Item} \mid \text{List1} ; \text{List1}$$

by

$$\text{List1} \rightarrow \text{Item Rest} \quad \text{Rest} \rightarrow \epsilon \mid ; \text{Item Rest}$$

Resolving ambiguity with added nonterminals

More serious example of ambiguity:

$$\begin{aligned} \text{Exp} &\rightarrow \text{Num} \mid \text{Var} \mid (\text{Exp}) \mid - \text{Exp} \\ &\mid \text{Exp} + \text{Exp} \mid \text{Exp} - \text{Exp} \mid \text{Exp} * \text{Exp} \mid \text{Exp}/\text{Exp} \end{aligned}$$

We can disambiguate this by adding nonterminals to capture more subtle distinctions between different classes of expressions:

$$\begin{aligned} \text{Exp} &\rightarrow \text{ExpA} \mid \text{Exp} + \text{ExpA} \mid \text{Exp} - \text{ExpA} \\ \text{ExpA} &\rightarrow \text{ExpB} \mid \text{ExpA} * \text{ExpB} \mid \text{ExpA} / \text{ExpB} \\ \text{ExpB} &\rightarrow \text{ExpC} \mid - \text{ExpC} \\ \text{ExpC} &\rightarrow \text{Num} \mid \text{Var} \mid (\text{Exp}) \end{aligned}$$

Note that this builds in certain **design decisions** concerning what we want the rules of precedence to be — shouldn't entrust this process to a machine!

N.B. our revised grammar is **unambiguous**, but not yet **LL(1)** ...

Problem 2: Shared prefixes

Consider the two productions

$$\text{Stmt} \rightarrow \text{do Stmt while Cond}$$

$$\text{Stmt} \rightarrow \text{do Stmt until Cond}$$

On encountering the nonterminal Stmt and the terminal do, an LL(1) parser would have no way of choosing between these two rules.

Solution: factor out the common part of these rules, so 'delaying' the decision until the relevant information becomes available:

$$\text{Stmt} \rightarrow \text{do Stmt Test}$$

$$\text{Test} \rightarrow \text{while Cond} \mid \text{until Cond}$$

This simple trick is known as **left factoring**.

Problem 3: Left recursion

Suppose our grammar contains a rule like

$$\text{Exp} \rightarrow \text{Exp} + \text{ExpA}$$

Problem: whatever terminals Exp could begin with, $\text{Exp} + \text{ExpA}$ could also begin with. So there's a danger our parser would apply this rule indefinitely:

$$\text{Exp} \Rightarrow \text{Exp} + \text{ExpA} \Rightarrow \text{Exp} + \text{ExpA} + \text{ExpA} \Rightarrow \dots$$

(In practice, we wouldn't even get this far: there'd be a clash in the parse table, e.g. at Num , Exp .)

So **left recursion** makes a grammar non-LL(1).

Eliminating left recursion

Consider e.g. the rules

$$\text{Exp} \rightarrow \text{ExpA} \mid \text{Exp} + \text{ExpA} \mid \text{Exp} - \text{ExpA}$$

Taken together, these say that Exp can consist of ExpA followed by zero or more suffixes $+\text{ExpA}$ or $-\text{ExpA}$.

So we just need to formalize this!

$$\text{Exp} \rightarrow \text{ExpA OpsA} \quad \text{OpsA} \rightarrow \epsilon \mid +\text{ExpA OpsA} \mid -\text{ExpA OpsA}$$

(Reminiscent of [Arden's rule](#).) Likewise:

$$\text{ExpA} \rightarrow \text{ExpB OpsB} \quad \text{OpsB} \rightarrow \epsilon \mid +\text{ExpB OpsB} \mid -\text{ExpB OpsB}$$

Together with the earlier rules for ExpB and ExpC , these give an **LL(1)** version of our grammar for arithmetic expressions.

Indirect left recursion

Left recursion can also arise in a more indirect way. E.g.

$$A \rightarrow a \mid Bc \qquad B \rightarrow b \mid Ad$$

By considering the combined effect of these rules, can see that they are equivalent to the following LL(1) grammar.

$$\begin{array}{ll} A \rightarrow aE \mid bcE & B \rightarrow bF \mid adF \\ E \rightarrow \epsilon \mid dcE & F \rightarrow \epsilon \mid cdF \end{array}$$

(Won't go into the systematic method here.)

LL(1) grammars: summary

- Often (not always), a “natural” grammar for some language of interest can be massaged into an LL(1) grammar. This allows for very efficient parsing.
- Knowing a grammar is LL(1) also assures us that it is **unambiguous** — often non-trivial! By the same token, LL(1) grammars are poorly suited to **natural languages**.
- However, an LL(1) grammar may be less readable and intuitive than the original. It may also appear to mutilate the ‘natural’ structure of phrases. We must take care not to mutilate it so much that we can no longer ‘execute’ the phrase as intended.
- One can design realistic computer languages with LL(1) grammars. For less cumbersome syntax that ‘flows’ better, one might want to go a bit beyond LL(1) (e.g. to LR(1)), but the principles remain the same.

Example of an LL(1) grammar

Here is a minor modification of the programming language grammar from Lecture 8. Combining it with our revised grammar for arithmetic expressions, we get an LL(1) grammar for a respectable programming language.

```

stmt   →  if-stmt | while-stmt | begin-stmt | assg-stmt
if-stmt →  if bool-expr then stmt else stmt
while-stmt →  while bool-expr do stmt
begin-stmt →  begin stmt-list end
stmt-list →  stmt stmts
stmts   →  ε | ; stmt stmts
assg-stmt →  VAR := arith-expr
bool-expr →  arith-expr compare-op arith-expr
compare-op →  < | > | <= | >= | == | != =

```

Clicker Question

Consider the alphabet of ASCII characters. Let N be the lexical class of all non-alphabetic characters. Consider the following context-free grammar for a nonterminal P .

$$P \rightarrow \epsilon \mid N P \mid P N$$

$$P \rightarrow a \mid a P a \mid a P A \mid A P a \mid A P A \mid A$$

$$P \rightarrow b \mid b P b \mid b P B \mid B P b \mid B P B \mid B$$

... (23 similar lines for 'C' to 'Y')

$$P \rightarrow z \mid z P z \mid z P Z \mid Z P z \mid Z P Z \mid Z$$

Which (if any) of the following ASCII strings cannot be parsed as a P ?

- 1 never odd or even
- 2 "Norma is as selfless as I am, Ron."
- 3 Live dirt up a side-track carted is a putrid evil.
- 4 I made reviled tubs repel; no, it is opposition, lepers, but delivered am I.
- 5 They can all be parsed.

Some light relief: Palindromic sentences

The grammar recognises palindromic alphabetic strings, ignoring whitespace, punctuation, case distinctions, etc.

It is not too hard to construct such strings consisting entirely of English words. However, it is rather satisfying to find examples that are coherent or interesting in some other way.

A famous example:

A man, a plan, a canal — Panama!

... which some smart aleck noticed could be tweaked to ...

A dog, a plan, a canal — Pagoda!

Probably there is nothing to equal ...

Best English palindrome in the world?

(From Guy Steele, *Common Lisp Reference Manual*, 1983.)

*A man, a plan, a canoe, pasta, heros, rajahs, a
coloratura, maps, snipe, percale, macaroni, a gag, a
banana bag, a tan, a tag, a banana bag again (or a
camel), a crepe, pins, Spam, a rut, a Rolo, cash, a jar,
sore hats, a peon, a canal — Panama!*

Clicker Question

Consider again our grammar for palindromic strings.

$$P \rightarrow \epsilon \mid NP \mid PN$$

$$P \rightarrow a \mid aPa \mid aPA \mid APa \mid APA \mid A$$

$$P \rightarrow b \mid bPb \mid bPB \mid BPb \mid BPB \mid B$$

... (23 similar lines for 'C' to 'Y')

$$P \rightarrow z \mid zPz \mid zPZ \mid ZPz \mid ZPZ \mid Z$$

Q. Is this grammar LL(1)?

- 1 Yes.
- 2 No.
- 3 Don't know.

Clicker Question

Consider again our grammar for palindromic strings.

$$P \rightarrow \epsilon \mid NP \mid PN$$

$$P \rightarrow a \mid aPa \mid aPA \mid APa \mid APA \mid A$$

$$P \rightarrow b \mid bPb \mid bPB \mid BPb \mid BPB \mid B$$

... (23 similar lines for 'C' to 'Y')

$$P \rightarrow z \mid zPz \mid zPZ \mid ZPz \mid ZPZ \mid Z$$

Q. Is this grammar LL(1)?

- 1 Yes.
- 2 No.
- 3 Don't know.

Q. Is it possible to provide an LL(1) grammar for the language of palindromes?

Addendum: Chomsky Normal Form

Whilst on the subject of 'transforming grammars into equivalent ones of some special kind' ...

A context-free grammar $\mathcal{G} = (N, \Sigma, P, S)$ is in **Chomsky normal form (CNF)** if all productions are of the form

$$A \rightarrow BC \quad \text{or} \quad A \rightarrow a \quad (A, B, C \in N, a \in \Sigma)$$

Theorem: Disregarding the empty string, every CFG \mathcal{G} is equivalent to a grammar \mathcal{G}' in Chomsky normal form. ($\mathcal{L}(\mathcal{G}') = \mathcal{L}(\mathcal{G}) - \{\epsilon\}$)

This is useful, because certain general parsing algorithms (e.g. the **CYK algorithm**, see Lecture 17) work best for grammars in CNF.

Converting to Chomsky Normal Form

Consider for example the grammar

$$S \rightarrow TT \mid [S] \quad T \rightarrow \epsilon \mid (T)$$

Step 1: remove all ϵ -productions, and for each rule $X \rightarrow \alpha Y \beta$, add a new rule $X \rightarrow \alpha \beta$ whenever Y 'can be empty'.

$$S \rightarrow TT \mid T \mid [S] \mid [] \quad T \rightarrow (T) \mid ()$$

Step 2: remove 'unit productions' $X \rightarrow Y$.

$$S \rightarrow TT \mid (T) \mid () \mid [S] \mid [] \quad T \rightarrow (T) \mid ()$$

Now all productions are of form $X \rightarrow a$ or $X \rightarrow x_1 \dots x_k$ ($k \geq 2$).

Converting to Chomsky Normal Form, ctd.

$$S \rightarrow TT \mid (T) \mid () \mid [S] \mid [] \quad T \rightarrow (T) \mid ()$$

Step 3: For each terminal a , add a nonterminal Z_a and a production $Z_a \rightarrow a$. In all rules $X \rightarrow x_1 \dots x_k$ ($k \geq 2$), replace each a by Z_a .

$$S \rightarrow TT \mid Z_l(TZ) \mid Z_l(Z) \mid Z_l[SZ_l] \mid Z_l[Z_l]$$

$$T \rightarrow Z_l(TZ) \mid Z_l(Z) \quad Z_l \rightarrow (\quad Z_l \rightarrow) \quad Z_l \rightarrow [\quad Z_l \rightarrow]$$

Step 4: For every production $X \rightarrow Y_1 \dots Y_n$ with $n \geq 3$, add new symbols W_2, \dots, W_{n-1} and replace the production with $X \rightarrow Y_1 W_2, \quad W_2 \rightarrow Y_2 W_3, \quad \dots, \quad W_{n-1} \rightarrow Y_{n-1} Y_n$.

E.g. $S \rightarrow Z_l(TZ) \mid Z_l[SZ_l]$ become

$$S \rightarrow Z_l W \quad W \rightarrow TZ) \quad S \rightarrow Z_l V \quad V \rightarrow SZ_l$$

The resulting grammar is now in **Chomsky Normal Form**.

Reading

- **Making grammars LL(1):** former lecture notes available via the Course Schedule webpage.
- **Chomsky Normal Form:** Kozen chapter 21, Jurafsky & Martin section 12.5.