

Agenda: an algorithm for recovering the analyses of an input sentence with respect to a given context-free grammar (CFG).

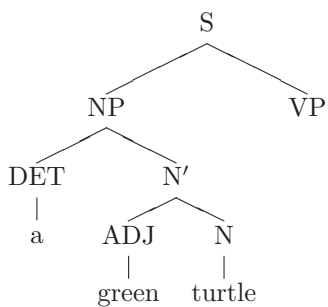
Announcements:

- The second prelim will be held in class on Friday, November 4th; topic coverage: IR and Web-based retrieval. Same drill as before (including the fact that no homework will be assigned before then), except that I don't anticipate being able to provide additional sample problems. Further details, including changes to the usual office hour schedule, are forthcoming.
- If you are either in favor or against changing the Monday 8pm regular office hour to Monday 1:30-2:30, please let me know by Thursday midnight by email (stating the reasons for your preference, e.g., conflict with the changed time but not the original time, or vice versa).

I. Garden-path sentences

1. "The cotton clothing is made of grows in Mississippi."
2. "The player kicked the ball kicked the ball."

II. Example Earley-style partial parse trees The sentence being parsed is "a green turtle swam to shore".



The point is that the leftmost leaves correspond to the first words of the sentence being parsed: Earley's algorithm tries to "guess" at how the parse tree will grow, but "anchors" its guesses against the words of the sentence in a left-to-right manner.

III. Parse states Assume a fixed sentence $w_1w_2 \dots w_n$ and CFG with start non-terminal S. The general form of a parse state is

$$(X \rightarrow \alpha \bullet \beta, i, j),$$

where

- α and β are some "stuff" (a sequence of zero or more terminals or nonterminals) such that $X \rightarrow \alpha\beta$ is a rewrite rule in the CFG,
- i and j range between 0 and n , and are *usually* interpreted as indicated endpoints of some subsequence of the sentence, and
- if $1 \leq i \leq j$, then we have inferred through the parsing process that α (i.e., the "stuff" before the "dot") can be rewritten into the sentence subsequence $w_iw_{i+1} \dots w_j$.

IV. Parser actions As above, we use Greek lower case letters to indicate “stuff” (a sequence of zero or more terminals or nonterminals). So, if we have rewrite rule $A \rightarrow BcDD$, we are allowed to describe it as $A \rightarrow \alpha B\beta$ or $A \rightarrow \alpha BcDD$ (the case in which α is a sequence of zero symbols) or a multitude of other ways.

1. *Predict* (guess how to expand the leftmost nonterminal not yet known to account for any part of the sentence):

Given $(X \rightarrow \alpha \bullet Y\beta, i, j)$ and rewrite rule $Y \rightarrow \gamma$, add state $(Y \rightarrow \bullet\gamma, j + 1, j)$.

Note the special treatment of what are usually endpoint indicators, meant to remind us that we don't know how far into the sentence that γ will “cover”.

2. *Scan* (advance the “dot” over a terminal that matches the next word in the sentence):

Given $(X \rightarrow \alpha \bullet y\beta, i, j)$ and the fact that $y = w_{j+1}$, add state $(X \rightarrow \alpha y \bullet \beta, i, j + 1)$.

3. *Complete* (advance the “dot” over a nonterminal that we now know can be rewritten into the next subsequence of words in the sentence):

Given $(X \rightarrow \alpha \bullet Y\beta, i, j)$ and $(Y \rightarrow \gamma \bullet, j + 1, k)$, add state $(X \rightarrow \alpha Y \bullet \beta, i, k)$.

V. Earley's algorithm This algorithm is guaranteed to terminate.

1. Start with the special state $(\rightarrow \bullet S, 1, 0)$.
2. Perform all possible predictions
3. Scan.
4. Perform all possible completions.
5. Return to step 2 unless nothing changed from the previous round.

VI. Example partial execution Note that the process would be more interesting for grammars generating ambiguous output.

For brevity (don't do this at home!), we just give the rewrite rules for a sample CFG with start nonterminal S:

- | | |
|-----------------------------|-----------------------------|
| (1) $S \rightarrow NP VP$ | (2) $ADJ \rightarrow green$ |
| (3) $NP \rightarrow DET N'$ | (4) $N \rightarrow turtle$ |
| (5) $N' \rightarrow ADJ N$ | (6) $V \rightarrow swam$ |
| (7) $VP \rightarrow V P N$ | (8) $P \rightarrow to$ |
| (9) $DET \rightarrow a$ | (10) $N \rightarrow shore$ |

Sentence: “a green turtle swam to shore”.

<i>init:</i> $(\rightarrow \bullet S, 1, 0)$	<i>predict:</i> $(N' \rightarrow \bullet ADJ N, 2, 1)$	<i>predict:</i> $(N \rightarrow \bullet turtle, 3, 2)$
<i>predict:</i> $(S \rightarrow \bullet NP VP, 1, 0)$	$(ADJ \rightarrow \bullet green, 2, 1)$	$(N \rightarrow \bullet shore, 3, 2)$
$(NP \rightarrow \bullet DET N', 1, 0)$	<i>scan:</i> $(ADJ \rightarrow green \bullet, 2, 2)$	<i>scan:</i> $(N \rightarrow turtle \bullet, 3, 3)$
$(DET \rightarrow \bullet a, 1, 0)$	<i>complete:</i> $(N' \rightarrow ADJ \bullet N, 2, 2)$	<i>complete:</i> $(N' \rightarrow ADJ N \bullet, 2, 3)$
<i>scan:</i> $(DET \rightarrow a \bullet, 1, 1)$		$(NP \rightarrow DET N' \bullet, 1, 3)$
<i>complete:</i> $(NP \rightarrow DET \bullet N', 1, 1)$		$(S \rightarrow NP \bullet VP, 1, 3)$

This table implicitly encodes both the partial parse tree given on the reverse of this sheet and an alternative that predicts that “shore” would be the third word, but this alternative is never expanded further.