---

**function** CYK(*words,grammar*) **returns** The most probable parse
                                                                           and its probability

Create and clear $\pi[num\_words, num\_words, num\_nonterminals]$

\# base case
**for** $i \leftarrow 1$ **to** *num_words*
   **for** $A \leftarrow 1$ **to** *num_nonterminals*
     **if** $(A \rightarrow w_i)$ is in grammar  **then**
       $\pi[i, i, A] \leftarrow P(A \rightarrow w_i)$

\# recursive case
**for** $span \leftarrow 2$ **to** *num_words*
   **for** $begin \leftarrow 1$ **to** $num\_words - span + 1$
    $end \leftarrow begin + span - 1$
    **for** $m =$ begin **to** $end - 1$
      **for** $A = 1$ **to** *num_nonterminals*
      **for** $B = 1$ **to** *num_nonterminals*
      **for** $C = 1$ **to** *num_nonterminals*
        $prob = \pi[begin, m, B] \times \pi[m+1, end, C] \times P(A \rightarrow B\,C)$
       **if** $(prob > \pi[begin, end, A])$ **then**
        $\pi[begin, end, A] = prob$
        $back[begin, end, A] = \{m, B, C\}$
**return** $build\_tree(back[1, num\_words, 1]), \pi[1, num\_words, 1]$

---

**Figure 12.3**    The Probabilistic CYK algorithm for finding the maximum probability parse of a string of *num_words* words given a PCFG grammar with *num_rules* rules in Chomsky Normal Form (after Collins (1999) and Aho and Ullman (1972).)  *back* is an array of back-pointers used to recover the best parse. The *build_tree* function is left as an exercise to the reader.

When a treebank is unavailable, the counts needed for computing PCFG probabilities can be generated by first parsing a corpus. If sentences were unambiguous, it would be as simple as this: parse the corpus, increment a counter for every rule in the parse, and then normalize to get probabilities. However, since most sentences are ambiguous, in practice we need to keep a separate count for each parse of a sentence and weight each partial count by the probability of the parse it appears in. The standard algorithm for computing this is called the **Inside-Outside** algorithm, and was proposed   INSIDE-OUTSIDE by Baker (1979) as a generalization of the forward-backward algorithm of Chapter 7. See Manning and Schütze (1999) for a complete description of