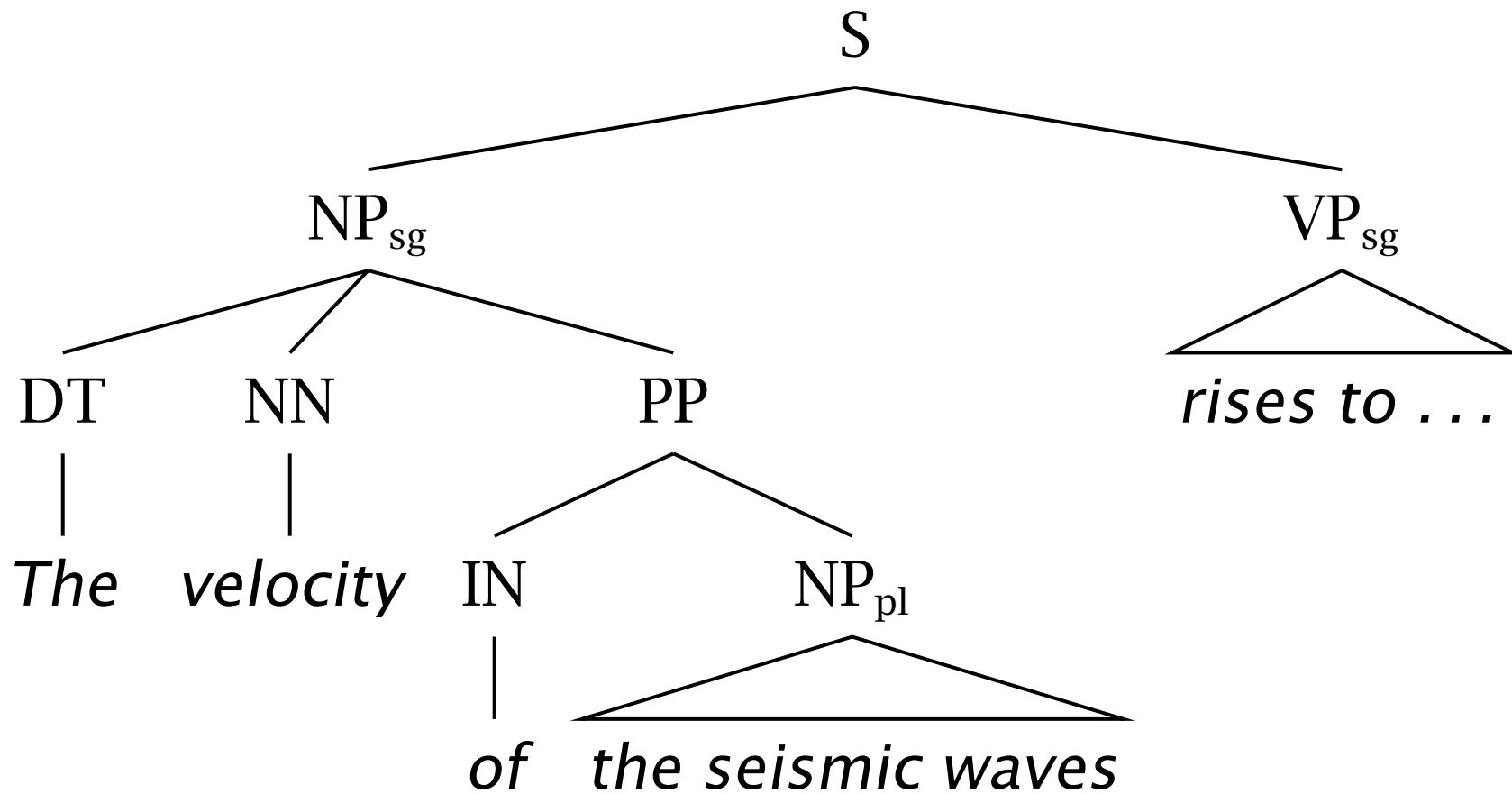


{Probabilistic|Stochastic} Context-Free Grammars (PCFGs)

- The velocity of the seismic waves rises to ...



PCFGs

A PCFG G consists of:

- A set of terminals, $\{w^k\}, k = 1, \dots, V$
- A set of nonterminals, $\{N^i\}, i = 1, \dots, n$
- A designated start symbol, N^1
- A set of rules, $\{N^i \rightarrow \zeta^j\}$, (where ζ^j is a sequence of terminals and nonterminals)
- A corresponding set of probabilities on rules such that:

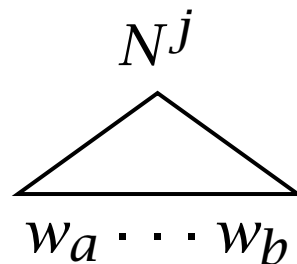
$$\forall i \quad \sum_j P(N^i \rightarrow \zeta^j) = 1$$

PCFG notation

Sentence: sequence of words $w_1 \cdots w_m$

w_{ab} : the subsequence $w_a \cdots w_b$

N_{ab}^i : nonterminal N^i dominates $w_a \cdots w_b$



$N^i \xRightarrow{*} \zeta$: Repeated derivation from N^i gives ζ .

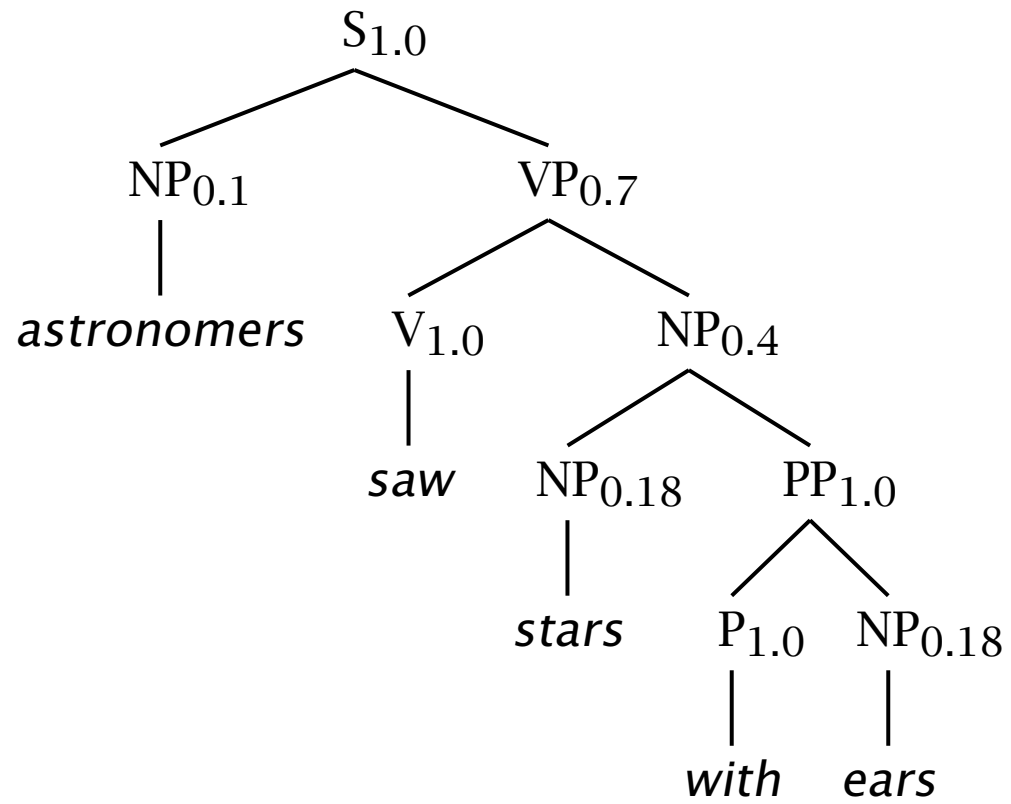
PCFG probability of a string

$$\begin{aligned} P(w_{1n}) &= \sum_t P(w_{1n}, t) \quad t \text{ a parse of } w_{1n} \\ &= \sum_{\{t:\text{yield}(t)=w_{1n}\}} P(t) \end{aligned}$$

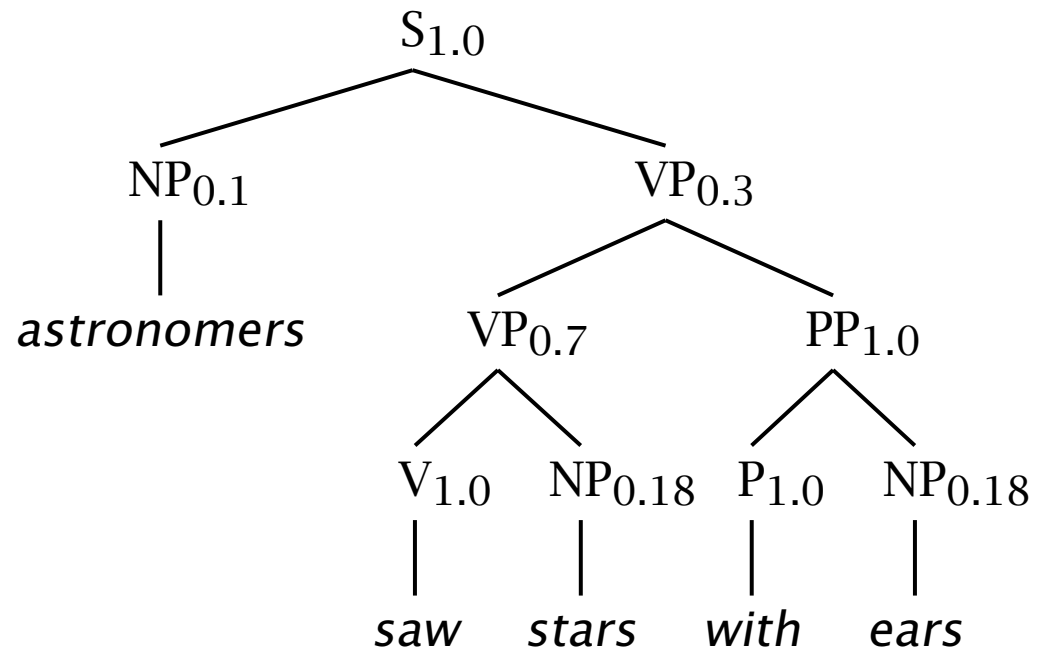
A simple PCFG (in CNF)

$S \rightarrow NP VP$	1.0	$NP \rightarrow NP PP$	0.4
$PP \rightarrow P NP$	1.0	$NP \rightarrow \textit{astronomers}$	0.1
$VP \rightarrow V NP$	0.7	$NP \rightarrow \textit{ears}$	0.18
$VP \rightarrow VP PP$	0.3	$NP \rightarrow \textit{saw}$	0.04
$P \rightarrow \textit{with}$	1.0	$NP \rightarrow \textit{stars}$	0.18
$V \rightarrow \textit{saw}$	1.0	$NP \rightarrow \textit{telescopes}$	0.1

t_1 :



*t*₂:



The two parse trees' probabilities and the sentence probability

$$\begin{aligned}P(t_1) &= 1.0 \times 0.1 \times 0.7 \times 1.0 \times 0.4 \\ &\quad \times 0.18 \times 1.0 \times 1.0 \times 0.18 \\ &= 0.0009072\end{aligned}$$

$$\begin{aligned}P(t_2) &= 1.0 \times 0.1 \times 0.3 \times 0.7 \times 1.0 \\ &\quad \times 0.18 \times 1.0 \times 1.0 \times 0.18 \\ &= 0.0006804\end{aligned}$$

$$P(w_{15}) = P(t_1) + P(t_2) = 0.0015876$$

Assumptions of PCFGs

1. Place invariance (like time invariance in HMM):

$$\forall k \quad P(N_{k(k+c)}^j \rightarrow \zeta) \text{ is the same}$$

2. Context-free:

$$P(N_{kl}^j \rightarrow \zeta | \text{words outside } w_k \dots w_l) = P(N_{kl}^j \rightarrow \zeta)$$

3. Ancestor-free:

$$P(N_{kl}^j \rightarrow \zeta | \text{ancestor nodes of } N_{kl}^j) = P(N_{kl}^j \rightarrow \zeta)$$

Let the upper left index in ${}^iN^j$ be an arbitrary identifying index for a particular token of a nonterminal.

Then,

$$\begin{aligned}
 P \left(\begin{array}{c} {}^1S \\ \swarrow \quad \searrow \\ {}^2NP \quad {}^3VP \\ \swarrow \quad \searrow \quad | \\ the \quad man \quad snores \end{array} \right) &= P({}^1S_{13} \rightarrow {}^2NP_{12} {}^3VP_{33}, {}^2NP_{12} \rightarrow the_1 man_2, {}^3VP_{33} \rightarrow snores) \\
 &= \dots \\
 &= P(S \rightarrow NP VP)P(NP \rightarrow the man)P(VP \rightarrow snores)
 \end{aligned}$$

Some features of PCFGs

Reasons to use a PCFG, and some idea of their limitations:

- Partial solution for grammar ambiguity: a PCFG gives some idea of the plausibility of a sentence.
- But not a very good idea, as not lexicalized.
- Better for grammar induction (Gold 1967)
- Robustness. (Admit everything with low probability.)

Some features of PCFGs

- Gives a probabilistic language model for English.
- In practice, a PCFG is a worse language model for English than a trigram model.
- Can hope to combine the strengths of a PCFG and a trigram model.
- PCFG encodes certain biases, e.g., that smaller trees are normally more probable.

Improper (inconsistent) distributions

■ $S \rightarrow \text{rhubarb} \quad P = \frac{1}{3}$

$S \rightarrow S S \quad P = \frac{2}{3}$

■ $\text{rhubarb} \quad \frac{1}{3}$

$\text{rhubarb rhubarb} \quad \frac{2}{3} \times \frac{1}{3} \times \frac{1}{3} = \frac{2}{27}$

$\text{rhubarb rhubarb rhubarb} \quad \left(\frac{2}{3}\right)^2 \times \left(\frac{1}{3}\right)^3 \times 2 = \frac{8}{243}$

...

■ $P(\mathcal{L}) = \frac{1}{3} + \frac{2}{27} + \frac{8}{243} + \dots = \frac{1}{2}$

■ Improper/inconsistent distribution

■ Not a problem if you estimate from parsed treebank: Chi and Geman 1998).

Questions for PCFGs

Just as for HMMs, there are three basic questions we wish to answer:

- $P(w_{1m}|G)$
- $\arg \max_t P(t|w_{1m}, G)$
- Learning algorithm. Find G such that $P(w_{1m}|G)$ is maximized.

Chomsky Normal Form grammars

We'll do the case of Chomsky Normal Form grammars, which only have rules of the form:

$$N^i \rightarrow N^j N^k$$

$$N^i \rightarrow w^j$$

Any CFG can be represented by a weakly equivalent CFG in Chomsky Normal Form. It's straightforward to generalize the algorithm (recall chart parsing).

PCFG parameters

We'll do the case of Chomsky Normal Form grammars, which only have rules of the form:

$$N^i \rightarrow N^j N^k$$

$$N^i \rightarrow w^j$$

The parameters of a CNF PCFG are:

$P(N^j \rightarrow N^r N^s | G)$ A n^3 matrix of parameters

$P(N^j \rightarrow w^k | G)$ An nt matrix of parameters

For $j = 1, \dots, n$,

$$\sum_{r,s} P(N^j \rightarrow N^r N^s) + \sum_k P(N^j \rightarrow w^k) = 1$$

Probabilistic Regular Grammar:

$$N^i \rightarrow w^j N^k$$

$$N^i \rightarrow w^j$$

Start state, N^1

HMM:

$$\sum_{w_{1n}} P(w_{1n}) = 1 \quad \forall n$$

whereas in a PCFG or a PRG:

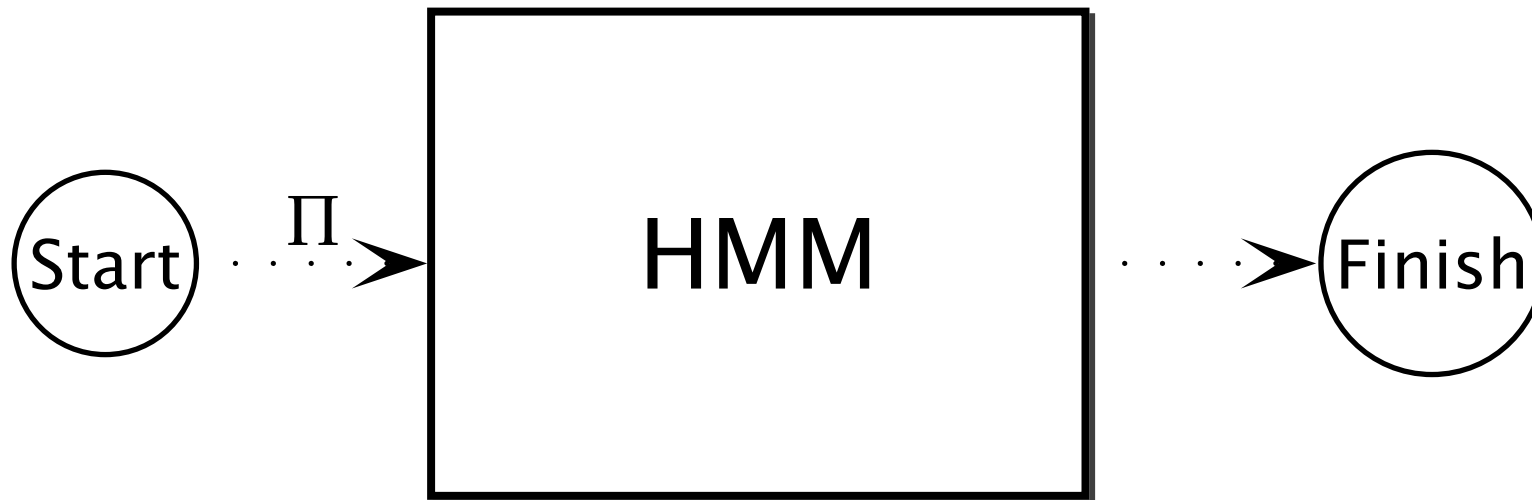
$$\sum_{w \in L} P(w) = 1$$

Consider:

$P(\text{John decided to bake a})$

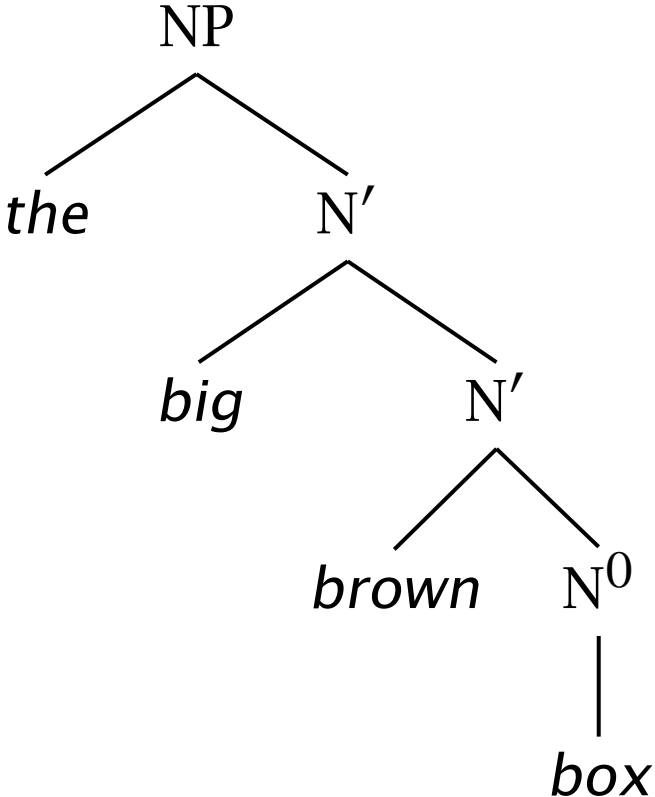
High probability in HMM, low probability in a PRG or a PCFG. Implement via sink state.

A PRG



Comparison of HMMs (PRGs) and PCFGs

$X: NP \rightarrow N' \rightarrow N' \rightarrow N' \rightarrow \text{sink}$
 $O: the \quad big \quad brown \quad box$



Inside and outside probabilities

This suggests: whereas for an HMM we have:

$$\text{Forwards} = \alpha_i(t) = P(w_{1(t-1)}, X_t = i)$$

$$\text{Backwards} = \beta_i(t) = P(w_{tT} | X_t = i)$$

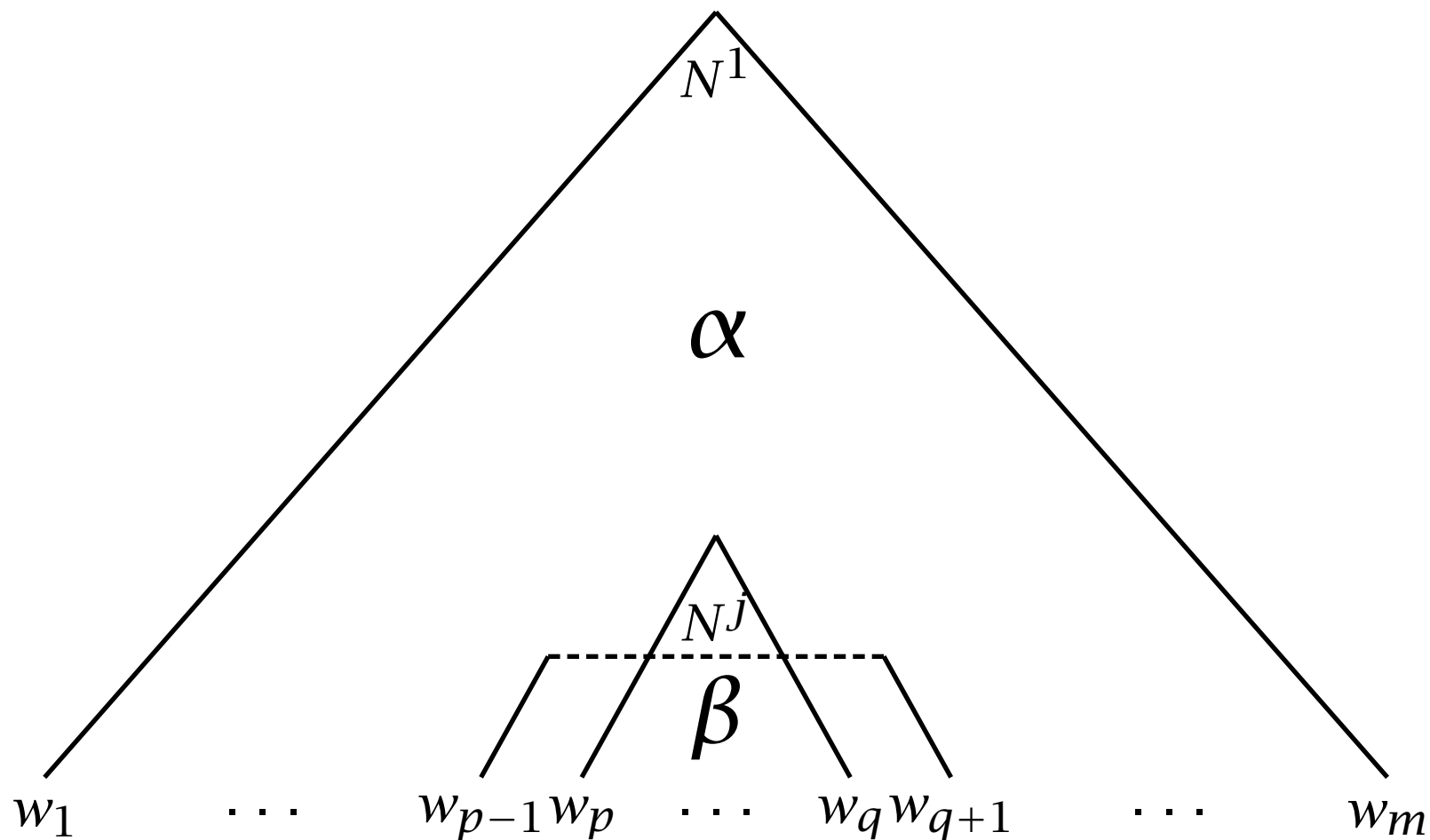
for a PCFG we make use of Inside and Outside probabilities, defined as follows:

$$\text{Outside} = \alpha_j(p, q) = P(w_{1(p-1)}, N_{pq}^j, w_{(q+1)m} | G)$$

$$\text{Inside} = \beta_j(p, q) = P(w_{pq} | N_{pq}^j, G)$$

A slight generalization of dynamic Bayes Nets covers PCFG inference by the inside-outside algorithm (and-or tree of conjunctive daughters disjunctively chosen)

Inside and outside probabilities in PCFGs.



Probability of a string

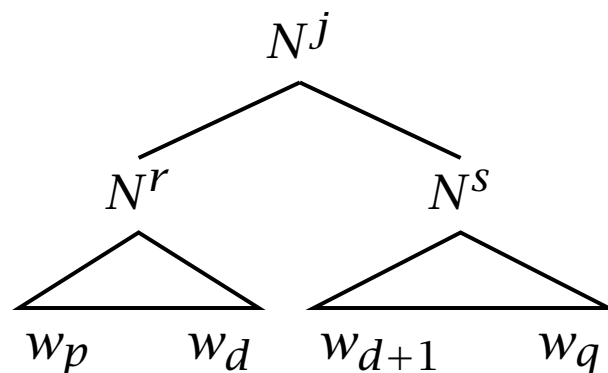
Inside probability

$$\begin{aligned} P(w_{1m}|G) &= P(N^1 \Rightarrow w_{1m}|G) \\ &= P(w_{1m}, N_{1m}^1, G) = \beta_1(1, m) \end{aligned}$$

Base case: We want to find $\beta_j(k, k)$ (the probability of a rule $N^j \rightarrow w_k$):

$$\begin{aligned} \beta_j(k, k) &= P(w_k | N_{kk}^j, G) \\ &= P(N^j \rightarrow w_k | G) \end{aligned}$$

Induction: We want to find $\beta_j(p, q)$, for $p < q$. As this is the inductive step using a Chomsky Normal Form grammar, the first rule must be of the form $N^j \rightarrow N^r N^s$, so we can proceed by induction, dividing the string in two in various places and summing the result:



These inside probabilities can be calculated bottom up.

For all j ,

$$\begin{aligned}
\beta_j(p, q) &= P(w_{pq} | N_{pq}^j, G) \\
&= \sum_{r,s} \sum_{d=p}^{q-1} P(w_{pd}, N_{pd}^r, w_{(d+1)q}, N_{(d+1)q}^s | N_{pq}^j, G) \\
&= \sum_{r,s} \sum_{d=p}^{q-1} P(N_{pd}^r, N_{(d+1)q}^s | N_{pq}^j, G) \\
&\quad P(w_{pd} | N_{pq}^j, N_{pd}^r, N_{(d+1)q}^s, G) \\
&\quad P(w_{(d+1)q} | N_{pq}^j, N_{pd}^r, N_{(d+1)q}^s, w_{pd}, G) \\
&= \sum_{r,s} \sum_{d=p}^{q-1} P(N_{pd}^r, N_{(d+1)q}^s | N_{pq}^j, G) \\
&\quad P(w_{pd} | N_{pd}^r, G) P(w_{(d+1)q} | N_{(d+1)q}^s, G) \\
&= \sum_{r,s} \sum_{d=p}^{q-1} P(N^j \rightarrow N^r N^s) \beta_r(p, d) \beta_s(d+1, q)
\end{aligned}$$

Calculation of inside probabilities (CKY algorithm)

	1	2	3	4	5
1	$\beta_{NP} = 0.1$		$\beta_S = 0.0126$		$\beta_S = 0.0015876$
2		$\beta_{NP} = 0.04$ $\beta_V = 1.0$	$\beta_{VP} = 0.126$		$\beta_{VP} = 0.015876$
3			$\beta_{NP} = 0.18$		$\beta_{NP} = 0.01296$
4				$\beta_P = 1.0$	$\beta_{PP} = 0.18$
5					$\beta_{NP} = 0.18$
	<i>astronomers</i>	<i>saw</i>	<i>stars</i>	<i>with</i>	<i>ears</i>

Outside probabilities

Probability of a string: For any k , $1 \leq k \leq m$,

$$\begin{aligned} P(w_{1m}|G) &= \sum_j P(w_{1(k-1)}, w_k, w_{(k+1)m}, N_{kk}^j | G) \\ &= \sum_j P(w_{1(k-1)}, N_{kk}^j, w_{(k+1)m} | G) \\ &\quad \times P(w_k | w_{1(k-1)}, N_{kk}^j, w_{(k+1)m}, G) \\ &= \sum_j \alpha_j(k, k) P(N^j \rightarrow w_k) \end{aligned}$$

Inductive (DP) calculation: One calculates the outside probabilities top down (after determining the inside probabilities).

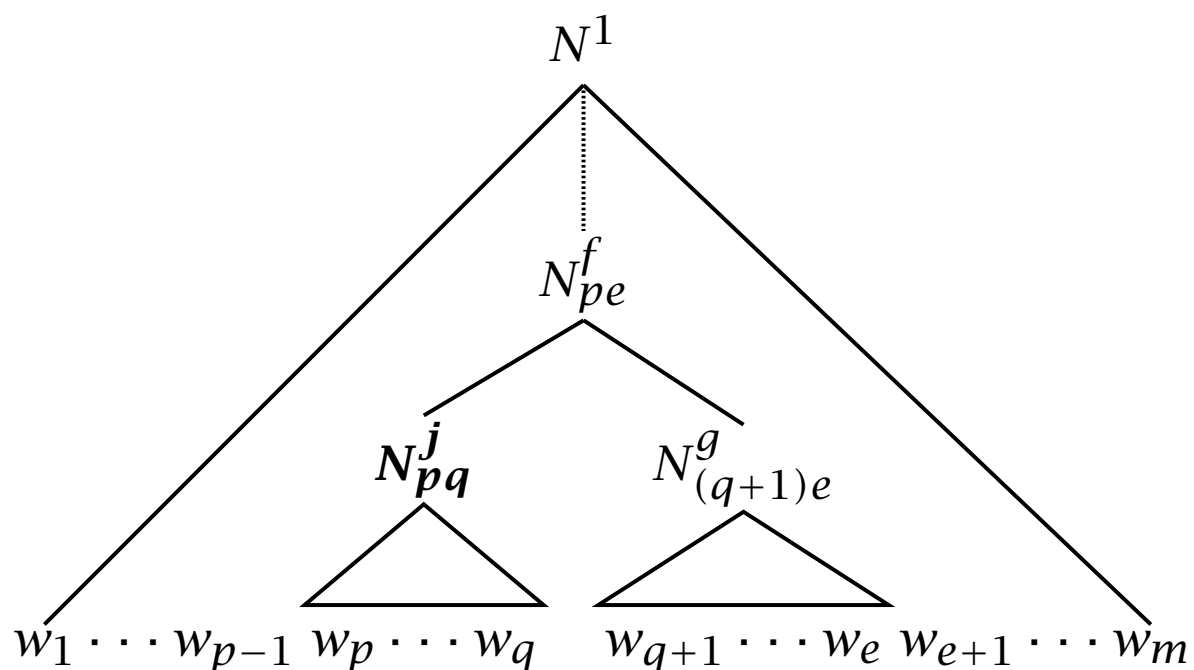
Outside probabilities

Base Case:

$$\alpha_1(1, m) = 1$$

$$\alpha_j(1, m) = 0, \text{ for } j \neq 1$$

Inductive Case:



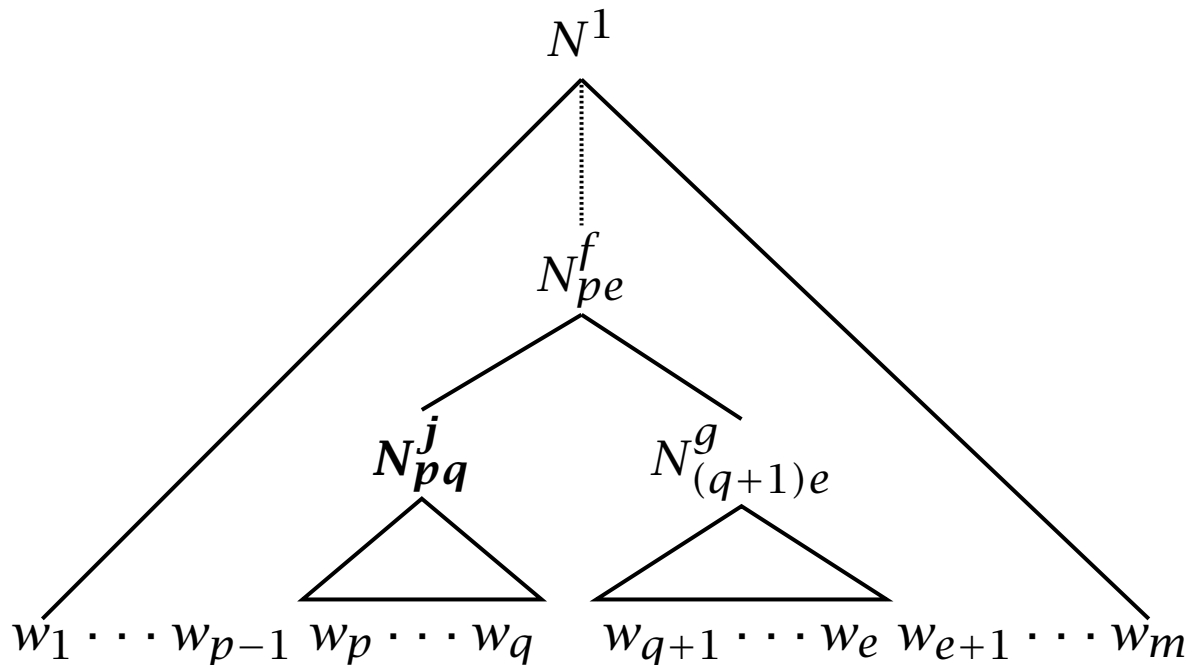
Outside probabilities

Base Case:

$$\alpha_1(1, m) = 1$$

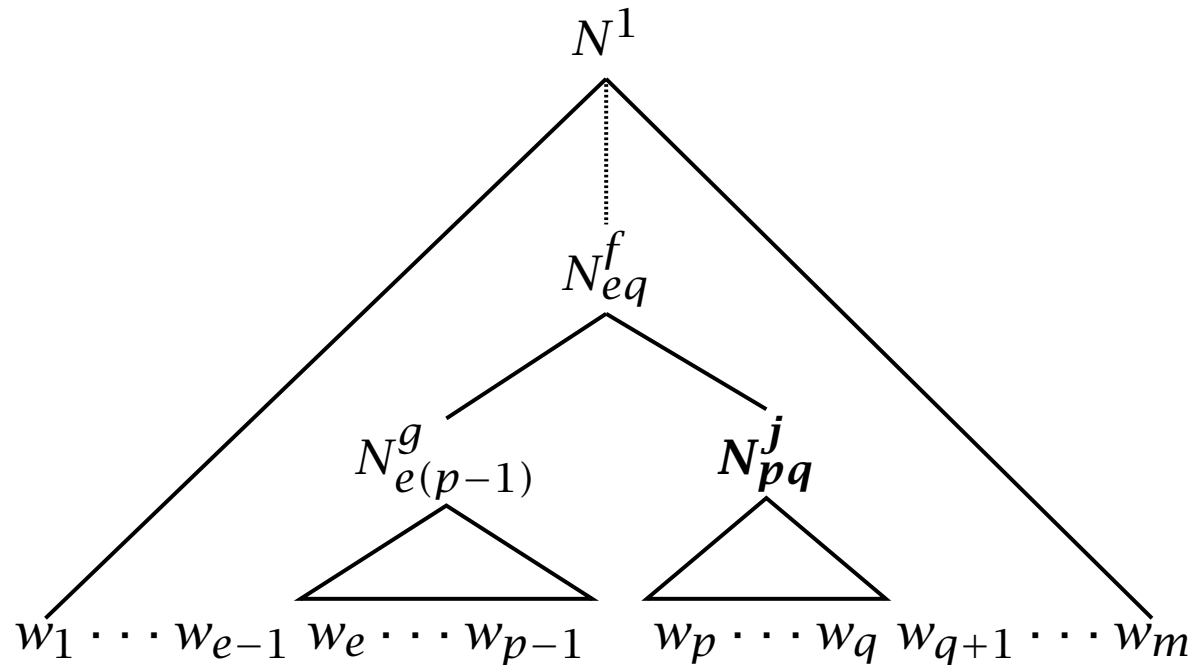
$$\alpha_j(1, m) = 0, \text{ for } j \neq 1$$

Inductive Case: it's either a left or right branch – we will solve over both possibilities and calculate using outside *and* inside probabilities



Outside probabilities – inductive case

A node N_{pq}^j might be the left or right branch of the parent node. We sum over both possibilities.



Inductive Case:

$$\begin{aligned}
 \alpha_j(p, q) &= \left[\sum_{f, g} \sum_{e=q+1}^m P(w_{1(p-1)}, w_{(q+1)m}, N_{pe}^f, N_{pq}^j, N_{(q+1)e}^g) \right. \\
 &\quad \left. + \left[\sum_{f, g} \sum_{e=1}^{p-1} P(w_{1(p-1)}, w_{(q+1)m}, N_{eq}^f, N_{e(p-1)}^g, N_{pq}^j) \right] \right] \\
 &= \left[\sum_{f, g, e} \sum_{j} \sum_{e=q+1}^m P(w_{1(p-1)}, w_{(e+1)m}, N_{pe}^f) P(N_{pq}^j, N_{(q+1)e}^g | N_{pe}^f) \right. \\
 &\quad \left. \times P(w_{(q+1)e} | N_{(q+1)e}^g) \right] + \left[\sum_{f, g} \sum_{e=1}^{p-1} P(w_{1(e-1)}, w_{(q+1)m}, N_{eq}^f) \right. \\
 &\quad \left. \times P(N_{e(p-1)}^g, N_{pq}^j | N_{eq}^f) P(w_{e(p-1)} | N_{e(p-1)}^g) \right] \\
 &= \left[\sum_{f, g} \sum_{e=q+1}^m \alpha_f(p, e) P(N^f \rightarrow N^j N^g) \beta_g(q+1, e) \right] \\
 &\quad + \left[\sum_{f, g} \sum_{e=1}^{p-1} \alpha_f(e, q) P(N^f \rightarrow N^g N^j) \beta_g(e, p-1) \right]
 \end{aligned}$$

Overall probability of a node existing

As with a HMM, we can form a product of the inside and outside probabilities. This time:

$$\begin{aligned} & \alpha_j(p, q) \beta_j(p, q) \\ &= P(w_{1(p-1)}, N_{pq}^j, w_{(q+1)m} | G) P(w_{pq} | N_{pq}^j, G) \\ &= P(w_{1m}, N_{pq}^j | G) \end{aligned}$$

Therefore,

$$p(w_{1m}, N_{pq} | G) = \sum_j \alpha_j(p, q) \beta_j(p, q)$$

Just in the cases of the root node and the preterminals, we know there will always be some such constituent.

Training a PCFG

We construct an EM training algorithm, as for HMMs. We would like to calculate how often each rule is used:

$$\hat{P}(N^j \rightarrow \zeta) = \frac{C(N^j \rightarrow \zeta)}{\sum_{\gamma} C(N^j \rightarrow \gamma)}$$

Have data \Rightarrow count; else work iteratively from expectations of current model.

Consider:

$$\begin{aligned}\alpha_j(p, q)\beta_j(p, q) &= P(N^1 \xRightarrow{*} w_{1m}, N^j \xRightarrow{*} w_{pq} | G) \\ &= P(N^1 \xRightarrow{*} w_{1m} | G)P(N^j \xRightarrow{*} w_{pq} | N^1 \xRightarrow{*} w_{1m}, G)\end{aligned}$$

We have already solved how to calculate $P(N^1 \Rightarrow w_{1m})$; let us call this probability π . Then:

$$P(N^j \xRightarrow{*} w_{pq} | N^1 \xRightarrow{*} w_{1m}, G) = \frac{\alpha_j(p, q)\beta_j(p, q)}{\pi}$$

and

$$E(N^j \text{ is used in the derivation}) = \sum_{p=1}^m \sum_{q=p}^m \frac{\alpha_j(p, q)\beta_j(p, q)}{\pi}$$

In the case where we are not dealing with a preterminal, we substitute the inductive definition of β , and $\forall r, s, p > q$:

$$P(N^j \rightarrow N^r N^s \Rightarrow w_{pq} | N^1 \Rightarrow w_{1n}, G) = \frac{\sum_{d=p}^{q-1} \alpha_j(p, q) P(N^j \rightarrow N^r N^s) \beta_r(p, d) \beta_s(d+1, q)}{\pi}$$

Therefore the expectation is:

$$E(N^j \rightarrow N^r N^s, N^j \text{ used})$$

$$\frac{\sum_{p=1}^{m-1} \sum_{q=p+1}^m \sum_{d=p}^{q-1} \alpha_j(p, q) P(N^j \rightarrow N^r N^s) \beta_r(p, d) \beta_s(d+1, q)}{\pi}$$

Now for the maximization step, we want:

$$P(N^j \rightarrow N^r N^s) = \frac{E(N^j \rightarrow N^r N^s, N^j \text{ used})}{E(N^j \text{ used})}$$

Therefore, the reestimation formula, $\hat{P}(N^j \rightarrow N^r N^s)$ is the quotient:

$$\hat{P}(N^j \rightarrow N^r N^s) =$$

$$\frac{\sum_{p=1}^{m-1} \sum_{q=p+1}^m \sum_{d=p}^{q-1} \alpha_j(p, q) P(N^j \rightarrow N^r N^s) \beta_r(p, d) \beta_s(d+1, q)}{\sum_{p=1}^m \sum_{q=1}^m \alpha_j(p, q) \beta_j(p, q)}$$

Similarly,

$$E(N^j \rightarrow w^k | N^1 \Rightarrow w_{1m}, G) =$$

$$\frac{\sum_{h=1}^m \alpha_j(h, h) P(N^j \rightarrow w_h, w_h = w^k)}{\pi}$$

Therefore,

$$\hat{P}(N^j \rightarrow w^k) = \frac{\sum_{h=1}^m \alpha_j(h, h) P(N^j \rightarrow w_h, w_h = w^k)}{\sum_{p=1}^m \sum_{q=1}^m \alpha_j(p, q) \beta_j(p, q)}$$

Inside-Outside algorithm: repeat this process until the estimated probability change is small.

Multiple training instances: if we have training sentences $W = (W_1, \dots, W_\omega)$, with $W_i = (w_1, \dots, w_{m_i})$ and we let u and v bet the common subterms from before:

$$u_i(p, q, j, r, s) = \frac{\sum_{d=p}^{q-1} \alpha_j(p, q) P(N^j \rightarrow N^r N^s) \beta_r(p, d) \beta_s(d+1, q)}{P(N^1 \Rightarrow W_i | G)}$$

and

$$v_i(p, q, j) = \frac{\alpha_j(p, q) \beta_j(p, q)}{P(N^1 \Rightarrow W_i | G)}$$

Assuming the observations are independent, we can sum contributions:

$$\hat{P}(N^j \rightarrow N^r N^s) = \frac{\sum_{i=1}^{\omega} \sum_{p=1}^{m_i-1} \sum_{q=p+1}^{m_i} u_i(p, q, j, r, s)}{\sum_{i=1}^{\omega} \sum_{p=1}^{m_i} \sum_{q=p}^{m_i} v_i(p, q, j)}$$

and

$$\hat{P}(N^j \rightarrow w^k) = \frac{\sum_{i=1}^{\omega} \sum_{\{h:w_h=w^k\}} v_i(h, h, j)}{\sum_{i=1}^{\omega} \sum_{p=1}^{m_i} \sum_{q=p}^{m_i} v_i(p, q, j)}$$

Problems with the Inside-Outside algorithm

1. Slow. Each iteration is $O(m^3 n^3)$, where $m = \sum_{i=1}^{\omega} m_i$, and n is the number of nonterminals in the grammar.
2. Local maxima are much more of a problem. Charniak reports that on each trial a different local maximum was found. Use simulated annealing? Restrict rules by initializing some parameters to zero? Or HMM initialization? Reallocate nonterminals away from “greedy” terminals?
3. Lari and Young suggest that you need many more nonterminals available than are theoretically necessary to get good grammar learning (about a threefold increase?). This compounds the first problem.
4. There is no guarantee that the nonterminals that the algorithm learns will have any satisfactory resemblance to the kinds of nonterminals normally motivated in linguistic analysis (NP, VP, etc.).