# Probabilistic Context Free Grammars

## 1  Defining PCFGs

A PCFG G consists of

1. A set of terminals: $\{w^k\}$, $k = 1 \ldots, V$

2. A set of non terminals: $\{N^i\}$, $i = 1 \ldots, n$

3. A designated Start symbol: $N^1$

4. A set of rules: $\{N^i \rightarrow \xi^j\}$, $i = 1 \ldots, n$

5. A probability function P which assigns probabilities to rules so that, for all nonterminals $N^i$

$$\sum_j P(N^i \rightarrow \xi^j) = 1$$

Conventions:

| Notation | Meaning |
|---|---|
| $G$ | Grammar (PCFG) |
| $\mathcal{L}$ | Language (generated or accepted by a grammar) |
| $t$ | parse tree |
| $\{N^1, \ldots, N^n\}$ | Nonterminal vocabulary ($N^1$ start symbol) |
| $\{w^1, \ldots, w^V\}$ | Terminal vocabulary |
| $w_1 \ldots w_m$ | Sentence to be parsed |
| $N_{pq}^j$ | Nonterminal $j$ spans positions $p$ through $q$ in sentence |
| $\alpha_j(p, q)$ | Outside probabilities |
| $\beta_j(p, q)$ | Inside probabilities |
| $N^j \overset{*}{\Rightarrow} w_p \ldots w_q$ | Nonterminal $j$ dominates words $w_p$ through $w_q$ in sentence (not necessarily directly) |

A key property of PCFGs is what we will call the **independence assumption**:

**Independence Assumption**
The probability of a node sequence $N^s\ N^r$ depends only on the immediate mother node, not any node above that or outside the current constituent.

We first address the problem of finding the probability of a string given a PCFG grammar. This is directly useful for some applications, such as filtering the hypoethses of a speech recognizer.

But considering some efficient algorithms for calculating string probabilities will have a side-effect. It will provide a line of attack on a somewhat more central problem for this course: finding the most probable parse of a string.

## 2 Example

Thus the only parameters of a PCFG grammar are the rule probabilities, as in the following simple PCFG:

(1)  
| | | | |
|---|---|---|---|
| $S \rightarrow NP\ VP$ | 1.0 | $NP \rightarrow NP\ PP$ | 0.4 |
| $PP \rightarrow P\ NP$ | 1.0 | $NP \rightarrow$ astronomers | 0.1 |
| $VP \rightarrow V\ NP$ | 0.7 | $NP \rightarrow$ ears | 0.18 |
| $VP \rightarrow VP\ PP$ | 0.3 | $NP \rightarrow$ saw | 0.04 |
| $P \rightarrow$ with | 1.0 | $NP \rightarrow$ stars | 0.18 |
| $V \rightarrow$ saw | 1.0 | $NP \rightarrow$ telescopes | 0.1 |

Notice the constraint

$$\sum_j P(NP \rightarrow \xi^j) = 1$$

is met. The sum of the rule probabilities for all the NP rules is 1.

We consider the probability of the sentence

(2)     Astronomers saw stars with ears.

with respect to this grammar, which admits the two analyses shown in Figures 1 and 2. The probability of each anaylsis is computed simply by multiplying the probabilities of the rules.

We compute the probability of the string by summing the probabilities of its two analyses:

$$
\begin{aligned}
P(w_{15}) &= P(t_1) + P(t_2) \\
&= 0.0009072 + 0.0006804 \\
&= 0.0015876
\end{aligned}
$$

## 3 Two algorithms for efficiently computing string probabilities

The naive approach to computing string probabilities: Compute the probability of each parse and add.

But: The number of parses of a string with a CFG grammar in general grows exponentially with the length of the string. Therefore this is inefficient.

### 3.1 Inside Probabilities

$$
\begin{aligned}
P(w_{1,m} \mid G) &= P(N^1 \overset{*}{\Rightarrow} w_{1,m} \mid G) \\
&= P(w_{1,m} \mid N^1_{1,m}, G) = \beta_1(1, m)
\end{aligned}
$$

The probability of a wordstring given the grammar equals the probability that the start symbol exhaustively dominates the string, which is the probability of the string $w_{1,m}$ given that the start symbol dominates words 1 through m, which we call the *inside probability* of the string for the span $1m$ and the category $N^1$. The general definition of an inside probability:

$$\beta_j(p, q) = P(w_{p,q} \mid N^j_{p,q}, G)$$

There is an an algotrithm for computing the inside probability of a strinng known as the **inside algorithm**. It breaks down to a base case and an induction:

1. **Base case**: The event that a nonterminal $N^j$ exhaustively dominates the word at position $k$ is $N^j_{kk}$ and the probability that that happens given the grammar is:

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

The point of the $\beta$ notation is that it uses only indices, indices of word positions (subscripts) and grammar elements (superscripts). The point of the 2nd line is that this reduces to a fact about the grammar, irrespective of word positions.

2. **Induction** We assume a Chomsky normal grammar, so the situation is as pictured in Figure 3.

   Then

$$
\beta_j(p,q) = P(w_{pq} \mid N^j_{pq}, G) \tag{1}
$$

$$
= \sum_{r,s}\sum_{d=p}^{q-1} P(w_{pd}, N^r_{pd}, w(d+1)q, N^s_{(d+1)q} \mid N^j_{pq}, G) \tag{2}
$$

We regard the probability of an analysis of $w_{pq}$ as an instance of category $N^j$ as the sum of the probabilities of an analysis using some way of expanding $N^j$; each such analysis uses some production:

$$
N^j \rightarrow N^r \, N^s
$$

Following the picture in Figure 3, we split the word string into two halves at word $d$, $p \leq d < q$, $w_{pd}$ and $w_{(d+1)q}$, one corresponding to $N^r$ and one to $N^s$. Thus the jointly occurring events in (2) are wordstring $w_{pd}$ being analyzed as $N^r$ and word string $w_{(d+1)q}$ being analyzed as $N^s$.

$$
\beta_j(p,q) = \sum_{r,s}\sum_{d=p}^{q-1} P(\ N^r_{pd},\ N^s_{(d+1)q} \mid N^j_{pq}, G) \times P(w_{pd} \mid N^j_{pq}, N^r_{pd}, N^s_{(d+1)q}, G) \\
\times P(w_{(d+1)q} \mid N^j_{pq}, N^r_{pd}, N^s_{(d+1)q}, w_{pd}, G) \tag{3}
$$

Equation (3) uses the chain rule on the joint events of (2). Now we use the independence assumptions of CFGs to conclude:

$$
\beta_j(p,q) = \sum_{r,s}\sum_{d=p}^{q-1} P(\ N^r_{pd},\ N^s_{(d+1)q} \mid N^j_{pq}, G) \times P(w_{pd} \mid N^r_{pd}, G) \\
\times P(w_{(d+1)q} \mid N^s_{(d+1)q}, G) \tag{4}
$$

And substituting in the definitions of an inside probability and a rule probability, we have the basic equation we can compute with:

$$
\beta_j(p,q) = \sum_{r,s}\sum_{d=p}^{q-1} P(\ N^j \rightarrow N^r \, N^s \mid N^j_{pq}, G) \times \beta_r(p,d) \times \beta_s(d+1,q) \tag{5}
$$

3

This means we can compute inside probabilities bottum up.

We use a dynamic programming algorithm, and as is customary represent the computations in a table. Cell $(i, j)$ for row $i$ and column $j$ contains the result of the $\beta$ computations for the string beginning with word $i$ and ending with word $j$. First we do all the word $\beta$s:

|   | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| 1 | $\beta_{NP} = \mathbf{0.1}$ | | | | |
| 2 | | $\beta_{NP} = \mathbf{0.04}$<br>$\beta_V = \mathbf{1.0}$ | | | |
| 3 | | | $\beta_{NP} = 0.18$ | | |
| 4 | | | | $\beta_P = 0.1$ | |
| 5 | | | | | $\beta_{NP} = 0.18$ |
| | astronomers | saw | stars | with | ears |

We now proceed as with CYK, trying to build constituents of length 2 next:

|   | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| 1 | $\beta_{NP} = 0.1$ | | | | |
| 2 | | $\beta_{NP} = 0.04$<br>$\beta_V = 1.0$ | $\beta_{VP} = \mathbf{0.126}$ | | |
| 3 | | | $\beta_{NP} = 0.18$ | | |
| 4 | | | | $\beta_P = 0.1$ | $\beta_{PP} = \mathbf{0.18}$ |
| 5 | | | | | $\beta_{NP} = 0.18$ |
| | astronomers | saw | stars | with | ears |

The computation of $\beta_{VP}(2,3)$ and $\beta_{PP}(4,5)$

$$
\begin{align}
\beta_{VP}(2,3) &= P(VP \to V\ NP) \times \beta_V(2,2) \times \beta_{NP}(3,3) \tag{6}\\
&= 0.7 \times 1.0 \times 0.18 \tag{7}\\
&= 0.126 \tag{8}\\
\beta_{PP}(4,5) &= P(PP \to P\ NP) \times \beta_P(4,4) \times \beta_{NP}(5,5) \tag{9}\\
&= 1.0 \times 1.0 \times 0.18 \tag{10}\\
&= 0.18 \tag{11}
\end{align}
$$

And in the next round:

|   | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| 1 | $\beta_{NP} = 0.1$ | | $\beta_S = \mathbf{0.0126}$ | | |
| 2 | | $\beta_{NP} = 0.04$<br>$\beta_V = 1.0$ | $\beta_{VP} = 0.126$ | | |
| 3 | | | $\beta_{NP} = 0.18$ | | $\beta_{NP} = \mathbf{0.01296}$ |
| 4 | | | | $\beta_P = 0.1$ | $\beta_{PP} = 0.18$ |
| 5 | | | | | $\beta_{NP} = 0.18$ |
| | astronomers | saw | stars | with | ears |

The computation of $\beta_S(1, 3)$ and $\beta_{NP}(3, 5)$

$$
\begin{aligned}
\beta_S(1, 3) &= P(S \rightarrow NP\ VP) \times \beta_{NP}(1, 1) \times \beta_{VP}(2, 3) \qquad &(12)\\
&= 1.0 \times 0.1 \times 0.126 &(13)\\
&= 0.0126 &(14)\\
\beta_{NP}(3, 5) &= P(NP \rightarrow NP\ PP) \times \beta_{NP}(3, 3) \times \beta_{PP}(4, 5) &(15)\\
&= 0.4 \times 0.18 \times 0.18 &(16)\\
&= 0.1296 &(17)
\end{aligned}
$$

And in the only round with an ambiguous constituent:

|   | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| 1 | $\beta_{NP} = 0.1$ | | $\beta_S = 0.0126$ | | |
| 2 | | $\beta_{NP} = 0.04$ $\beta_V = 1.0$ | $\beta_{VP} = 0.126$ | | $\boldsymbol{\beta_{VP} = 0.015876}$ |
| 3 | | | $\beta_{NP} = 0.18$ | | $\beta_{NP} = 0.01296$ |
| 4 | | | | $\beta_P = 0.1$ | $\beta_{PP} = 0.18$ |
| 5 | | | | | $\beta_{NP} = 0.18$ |
|   | astronomers | saw | stars | with | ears |

The computation of $\beta_{VP}(2, 5)$ :

$$
\begin{aligned}
\beta_{VP}(2, 5) &= (P(VP \rightarrow V\ NP) \times \beta_V(2, 2) \times \beta_{NP}(3, 5)) + \qquad &(18)\\
&\quad (P(VP \rightarrow VP\ PP) \times \beta_{VP}(2, 3) \times \beta_{PP}(4, 5))\\
&= (0.7 \times 1.0 \times 0.01296) + (0.3 * 0.126 * 0.18) &(19)\\
&= 0.009072 + 0.006804 &(20)\\
&= 015876 &(21)\\
& &(22)
\end{aligned}
$$

The **probability table** (or chart) may be completed in the next step:

|   | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| 1 | $\beta_{NP} = 0.1$ | | $\beta_S = 0.0126$ | | $\boldsymbol{\beta_S = 0.0015876}$ |
| 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 | | | | $\beta_P = 0.1$ | $\beta_{PP} = 0.18$ |
| 5 | | | | | $\beta_{NP} = 0.18$ |
|   | astronomers | saw | stars | with | ears |

The computation of $\beta_S(1, 5)$ is left as an exercise.

## 3.2 Outside Probabilities

We define $\alpha_j(p, q)$, the outside probability for a category $N^j$ and a span $w_{pq}$ as

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

This is the probability of generating $N^j$ covering the span $(p, q)$ (whatever the words in it) along with the all the words outside it.

Using the notion outside probability, we can view the problem of computing the probability of a string from the point of view of any word in it. Pick a word $w_k$. The following is true:

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

We start with the observation that any assignment of a category to $w_k$ fixes a probability for analyses of the entire string compatible with that assignment. So we sum up such analysis probabilities for each category $N^j$ in the grammar. We then apply the chain rule, and the independence assumptions of PCFGs, and the definition of $\alpha_j(p, q)$ to derive an expression for the probability of the entire string.

So if we knew the $\alpha_j(k, k)$ values, computing the string probability would be easy.

We now show the **outside algorithm**, a method for solving the more general problem of computing $\alpha_j(p, q)$. This is done topdown. We break it down as before into a base case and an induction:

1. **Base case**:
$$
\begin{aligned}
\alpha_1(1, m) &= 1 \\
&= 0 \text{ for } j \neq 1
\end{aligned}
$$

2. **Induction** The situation is as pictured in Figure 4 or Figure 5

   Then

$$
\begin{aligned}
\alpha_j(p, q) &= [\sum_{f,g \neq j} \sum_{e=q+1}^m P(w_{1(p-1)}, w_{(q+1)m}, N_{pe}^f, N_{pq}^j, N_{(q+1)e}^g)] \quad (23) \\
&\quad + [\sum_{f,g} \sum_{e=1}^{p-1} P(w_{1(p-1)}, w_{(q+1)m}, N_{eq}^f, N_{e(p-1)}^g, N_{pq}^j)]
\end{aligned}
$$

   The first double sum is the case of Figure 4. The second is the case of Figure 5. In the first case, we restrict the two categories not to be equal, so as not to count rules of the form:

$$X \to Y\,Y$$

   twice.

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

# 4 Finding the most probable parse

We define $\delta_j(p, q)$ as the highest probability parse of $w_{pq}$ as $N^j$

What follows is a Viterbi like algorithm for finding the best parse. Corresponding to Viterbi valkues we will have $\delta$ values. Corresponding to a Viterbi bnacktrace we will have poointers to the subtrees that help us produce our highest probvability parse for each (category,start, end) triple.

1. **Initialization**:
$$\delta_i(p, p) = P(N^i \rightarrow w_p)$$

2. **Induction**:
$$\delta_i(p, q) = \max_{\substack{1 \leq j,k \leq n \\ p \leq r < q}} P(N^i \rightarrow N^j \, N^k)\delta_j(p, r)\delta_k(r + 1, q)$$

Store backtrace:
$$\Psi_i(p, q) = \operatorname*{argmax}_{j,k,r} P(N^i \rightarrow N^j \, N^k)\delta_j(p, r)\delta_k(r + 1, q)$$
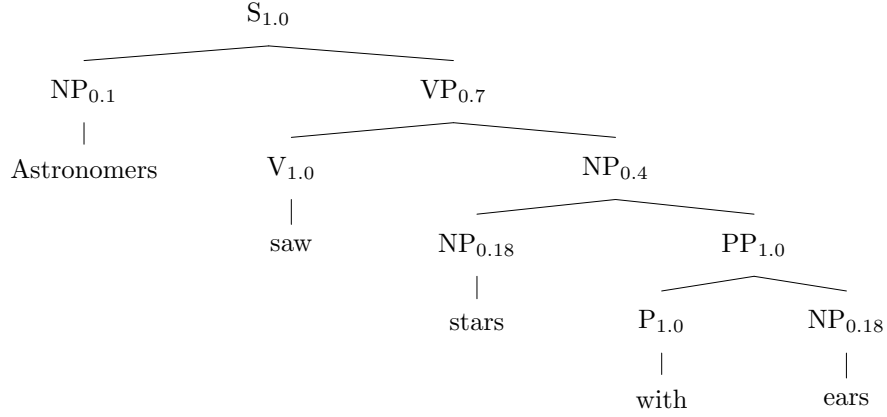
3. **Termination**: The probability of the max probability tree $\hat{t}$ is:
$$P(\hat{t}) = \delta_1(1, m)$$

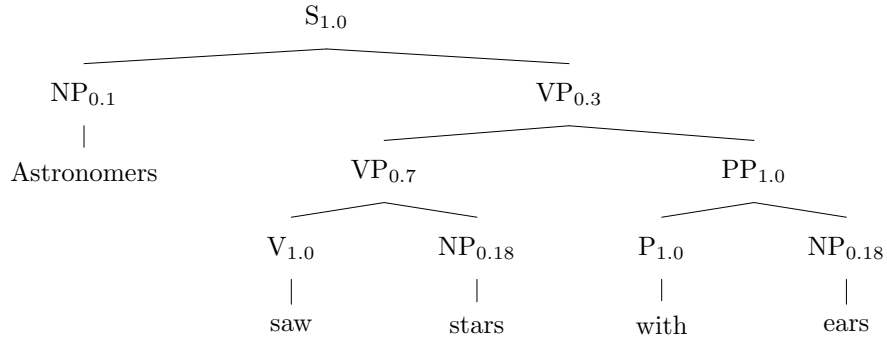In addition we know the mother node and span of the highest probability parse: $N^1$ and $(p, q)$.

In general given a mother node and a span $(i, p, q)$, we can construct a highest probability parse tree below it by following the $\Psi$ pointers:

$$
\begin{aligned}
&\text{Let} && \Psi_i(p, q) = (j, k, r) \\
&\text{Then} && \text{left daughter} = N^j_{pr} \\
& && \text{right daughter} = N^k_{(r+1)q}
\end{aligned}
$$

S$_{1.0}$

NP$_{0.1}$          VP$_{0.7}$

Astronomers     V$_{1.0}$              NP$_{0.4}$

saw        NP$_{0.18}$        PP$_{1.0}$

stars     P$_{1.0}$     NP$_{0.18}$

with        ears

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

Figure 1: $t_1$

S$_{1.0}$

NP$_{0.1}$                      VP$_{0.3}$

Astronomers       VP$_{0.7}$                PP$_{1.0}$

V$_{1.0}$      NP$_{0.18}$      P$_{1.0}$      NP$_{0.18}$

saw        stars       with        ears

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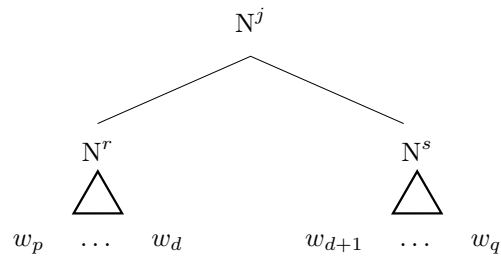&= 0.0006804
\end{aligned}
$$

Figure 2: $t_2$

N$^j$

N$^r$                N$^s$

$w_p$ ... $w_d$     $w_{d+1}$ ... $w_q$

Figure 3: Computation of inside probability: $\beta_j(p, q)$

$N^1$

$w_1 \ldots w_{p-1}$

$\vdots$

$N^f$

$w_{e+1} \ldots w_m$

$\mathbf{N^j}$

$\triangle$

$N^g$

$\triangle$

$w_p \quad \ldots \quad w_q$

$w_{q+1} \quad \ldots \quad w_e$

Figure 4: Computation of outside probability (left category): $\alpha_j(p, q)$

N$^1$

$w_1 \ldots w_{e-1}$ $\vdots$ $w_{q+1} \ldots w_m$

N$^f$

N$^g$ **N$^{\mathbf{j}}$**
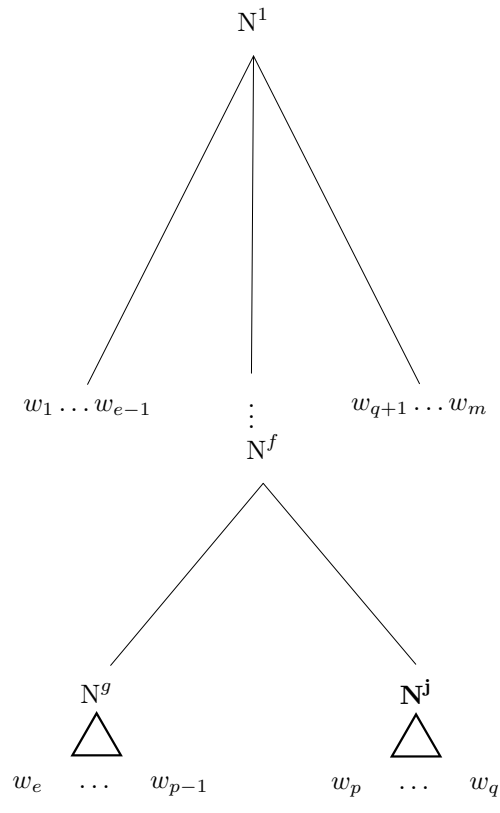
$w_e$ $\ldots$ $w_{p-1}$ $w_p$ $\ldots$ $w_q$

Figure 5: Computation of outside probability (right category): $\alpha_j(p, q)$