

## Chapter 6

# Syntax

Syntax, in linguistics, is the study of the structure of natural language sentences. It is sometimes approached as an “autonomous” part of language, describing what sentences are or aren’t grammatical, and sometimes approached as working together with semantics, describing the structure which is used to compute the meaning of a sentence.

## 6.1 Why Syntax?

Syntax proposes that there is something more to a sentence than just the sequence of words. To get an idea of what that something is, let’s look at some examples.

### 6.1.1 A quiz (not for credit)

For each pair of sentences, choose which one is better.

- I don’t like being made fun of.
  - I don’t like being taken pictures of.
- Who’d like some ice cream?
  - I’d!
  - I would!
- I saw a brown big spider.
  - I saw a big brown spider.
- That movie was fan-freaking-tastic.
  - That movie was fantas-freaking-tic.
- That’s the teacher that I couldn’t understand what he was talking about.
  - That’s the teacher that I couldn’t understand what was talking about.

### 6.1.2 Newspaper ambiguities

Consider the following example sentences, taken from newspaper articles.<sup>1</sup> What exactly is going wrong in these sentences, and how can we hope to make computers understand them?

- (6.1) Two cars were reported stolen by the Groveton police yesterday.
- (6.2) Mrs. Consigny was living alone in her home in Nakoma after her husband died in 1954 when the phone rang.
- (6.3) Black Panther leader Huey Newton, terming a 1974 murder charge “strictly a fabrication,” said yesterday he will testify at his trial on charges of killing a prostitute against his lawyer’s advice.
- (6.4) Yoko Ono will talk about her husband John Lennon who was killed in an interview with Barbara Walters.

### 6.1.3 Implementing subject-aux inversion

There are situations where we want computers to perform transformations on text that seem to require knowledge of concepts like “subject” or “object” – for example, machine translation. To keep things monolingual, imagine that we want a computer to turn statements into questions. How would you write a program to do this?

- (6.5)
  - a. The Pope is Catholic.
  - b. Is the Pope Catholic?
- (6.6)
  - a. Soylent Green is people.
  - b. Is Soylent Green people?
- (6.7)
  - a. Good fences make good neighbors.
  - b. Make good fences good neighbors?
  - c. Do good fences make good neighbors?
- (6.8)
  - a. A man who is his own lawyer has a fool for a client.
  - b. Is a man who his own lawyer has a fool for a client?
  - c. Does a man who is his own lawyer have a fool for a client?

### 6.1.4 Tests for constituency

To formulate a rule to do the above, we impose a tree structure on sentences, in which the leaves of the tree are words. (There’s an alternative kind of structure, called *dependency* trees, in which the internal nodes are also words.) The substring spanned by a node of the tree is called a *constituent*, and since we can’t see constituents, linguists have developed various *tests* for constituency:

Can you move it around?

<sup>1</sup><http://www.ling.upenn.edu/~beatrice/humor/newspaper-screwups.html>

(6.9) Fences make, good good neighbors.

(6.10) Make good, good fences neighbors.

(6.11) Good neighbors, good fences make.

(6.12) Fences make good, good neighbors.

(6.13) Make good neighbors, good fences.

Can you replace it with something like a pronoun?

(6.14) They make good neighbors.

(6.15) Good they good neighbors.

(6.16) Good fences it/they/. . . neighbors.

(6.17) Good fences make them.

(6.18) They good neighbors.

(6.19) Good it/they/. . . neighbors.

(6.20) Good fences do.

Can it participate in a *cleft*?

(6.21) It's good fences that make good neighbors.

(6.22) It's fences make that good good neighbors.

(6.23) It's make good that good fences neighbors.

(6.24) It's good neighbors that good fences make.

(6.25) It's good fences make that [do] good neighbors.

(6.26) It's make good neighbors that good fences do.

Can it be the answer to a question?

(6.27) What makes good neighbors? Good fences.

(6.28) Good what good neighbors? Fences make.

(6.29) Good fences what neighbors? Make good.

(6.30) Good fences make what? Good neighbors.

(6.31) What [do] good neighbors? Good fences make.

(6.32) Good fences do what? Make good neighbors.

## 6.2 Context Free Grammars

### 6.2.1 Why, what's wrong with finite automata?

We have spent a long time talking about all the things that finite automata can do, but there are important things that they can't do. For example, they cannot generate the language

$$L = \{a^n b^n \mid n \geq 0\} \quad (6.33)$$

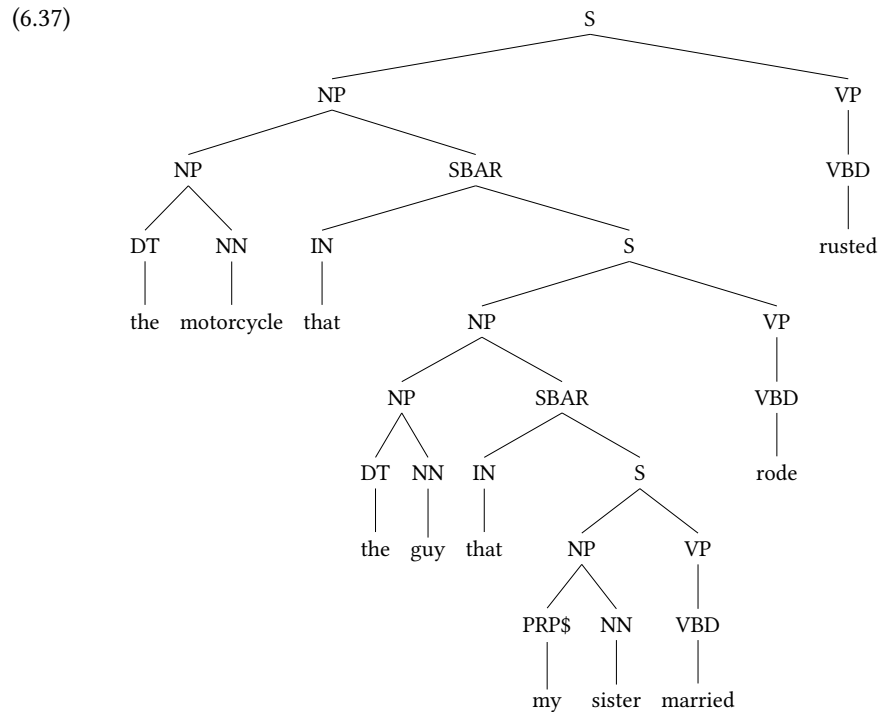
Linguistically, the analogous property that finite automata are missing is the ability to do *center-embedding*. English allows sentences like:

(6.34) The motorcycle rusted.

(6.35) The motorcycle that the guy rode rusted.

(6.36) The motorcycle that the guy that my sister married rode rusted.

At minimum, we need the number of noun phrases to equal the number of verbs, which we have already seen finite automata aren't able to do. Actually, we want to be able to create a structure like



which can tell us which noun corresponds to which verb. Later, we will see how this structure is also useful for translating into another language like Japanese or Hindi.

**Question 1.** The above argument holds only if you believe that *unbounded* center embedding is possible. In fact, center-embedding examples degrade rather quickly as more levels are added:

- (6.38) The cheese lay in the house that Jack built.
- (6.39) The cheese that the rat ate lay in the house that Jack built.
- (6.40) The cheese that the rat that the cat killed ate lay in the house that Jack built.
- (6.41) The cheese that the rat that the cat that the dog worried killed ate lay in the house that Jack built.
- ⋮
- (6.42) The cheese that the rat that the cat that the dog that the cow with the crumpled horn that the maiden all forlorn that the man all tattered and torn that the priest all shaven and shorn that the rooster that crowed in the morn that the farmer sowing his corn that the horse and the hound and the horn belonged to kept woke married kissed milked tossed worried killed ate lay in the house that Jack built.

If center embedding is bounded (say, to three levels), then how would you write a finite automaton to model it? What would still be unsatisfactory about such an account?

### 6.2.2 Context free grammars

Our solution is to use *context free grammars* (CFGs). CFGs are also widely used in compilers, where they are known as *Backus-Naur Form*. We begin with two examples of CFGs. The first one generates the non-finite-state language  $\{a^n b^n \mid n \geq 0\}$ :

$$\begin{aligned} S &\rightarrow aSb \\ S &\rightarrow \epsilon \end{aligned} \tag{6.43}$$

Here,  $S$  is called a *nonterminal symbol* and can be rewritten using one of the above rules, whereas  $a$  and  $b$  are called *terminal symbols* and cannot be rewritten. This is how the grammar works: start with a single  $S$ , then repeatedly choose the leftmost nonterminal and rewrite it using one of the rules until there are no more nonterminals. For example:

$$\begin{aligned} S &\Rightarrow aSb \\ &\Rightarrow aaSbb \\ &\Rightarrow aaaSbbb \\ &\Rightarrow aaabbbb \end{aligned} \tag{6.44}$$

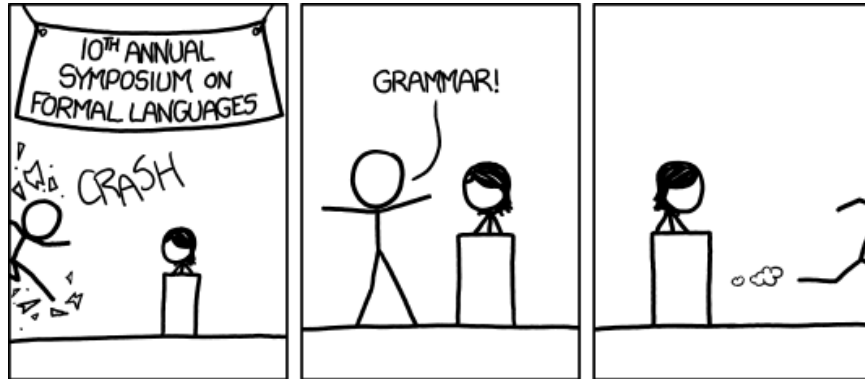


Figure 6.1: Formal Languages. [audience looks around] ‘What just happened?’  
‘There must be some context we’re missing.’

Our next example generates sentences (6.34), (6.35), and (6.36).

$$\begin{aligned}
 S &\rightarrow NP VP \\
 NP &\rightarrow DT NN \\
 NP &\rightarrow PRP\$ NN \\
 NP &\rightarrow NP SBAR \\
 VP &\rightarrow VBD \\
 SBAR &\rightarrow IN S \\
 DT &\rightarrow \text{the} \\
 PRP\$ &\rightarrow \text{my} \\
 NN &\rightarrow \text{motorcycle} \mid \text{guy} \mid \text{sister} \\
 VBD &\rightarrow \text{married} \mid \text{rode} \mid \text{rusted} \\
 IN &\rightarrow \text{that}
 \end{aligned} \tag{6.45}$$

Here, the uppercase symbols are nonterminal symbols, and the English words are terminal symbols. Also, we have used some shorthand:  $A \rightarrow \beta_1 \mid \beta_2$  stands for two rules,  $A \rightarrow \beta_1$  and  $A \rightarrow \beta_2$ .

**Question 2.** How would you use the above grammar to derive sentence (6.35)?

Here’s a more formal definition of CFGs.

**Definition 1.** A *context-free grammar* is a tuple  $G = (N, \Sigma, R, S)$ , where

- $N$  is a set of *nonterminal symbols*
- $\Sigma$  is a set of *terminal symbols*
- $R$  is a set of *rules* or *productions* of the form  $A \rightarrow \beta$ , where  $A \in N$  and  $\beta \in (N \cup \Sigma)^*$
- $S \in N$  is a distinguished *start symbol*

If  $A \in N$  and  $\alpha, \beta, \gamma \in (N \cup \Sigma)^*$ , we write  $\alpha A \gamma \Rightarrow_G \alpha \beta \gamma$  iff  $(A \rightarrow \beta) \in R$ , and we write  $\Rightarrow_G^*$  for the reflexive, transitive closure of  $\Rightarrow_G$ . Then the language generated by  $G$  is  $L(G) = \{w \in \Sigma^* \mid S \Rightarrow_G^* w\}$ .

### 6.2.3 Structure and ambiguity

As has already been alluded to, CFGs are interesting not only because they can generate more string languages than finite automata can, but because they build trees, known as *syntactic analyses*, *phrase-structure trees*, or *parse trees*. Whenever we use a rule  $A \rightarrow \beta$  to rewrite a nonterminal  $A$ , we don't erase  $A$  and replace it with  $\beta$ ; instead, we make the symbols of  $\beta$  the children of  $A$ .

**Question 3.** What would the tree for sentence (6.35) be?

One of the main purposes of these trees is that every subtree of the parse tree is supposed to have a *semantics* or meaning, so that the tree shows how to interpret the sentence. As a result, it is possible that a single string can have more than one structure, and therefore more than one meaning. This is called *ambiguity*. To illustrate it, we need a new example.

$$\begin{aligned}
 S &\rightarrow NP VP \\
 NP &\rightarrow DT NN \\
 NP &\rightarrow NN \\
 NP &\rightarrow NN NNS \\
 VP &\rightarrow VBP NP \\
 VP &\rightarrow VBP \\
 VP &\rightarrow VP PP \\
 PP &\rightarrow IN NP \\
 DT &\rightarrow a \mid an \\
 NN &\rightarrow time \mid fruit \mid arrow \mid banana \\
 NNS &\rightarrow flies \\
 VBP &\rightarrow flies \mid like \\
 IN &\rightarrow like
 \end{aligned}
 \tag{6.46}$$

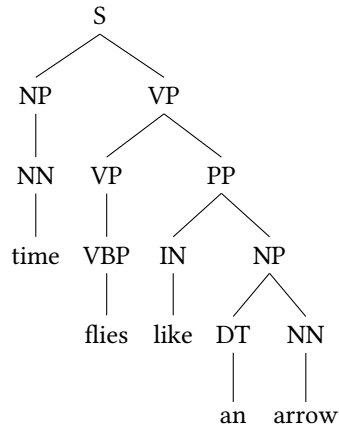
This grammar generates (among others) the following two strings:

(6.47) Time flies like an arrow.

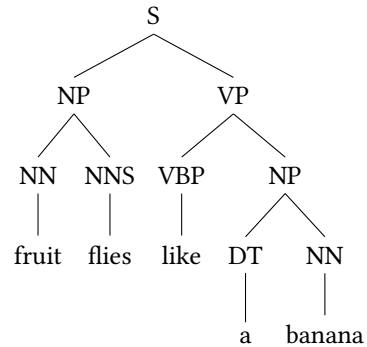
(6.48) Fruit flies like a banana.

Their “natural” structures are:

(6.49)

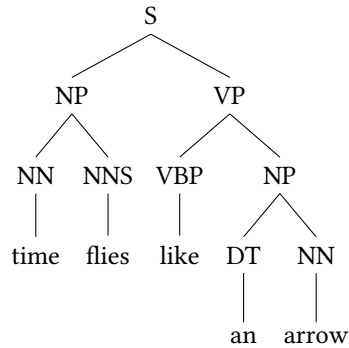


(6.50)

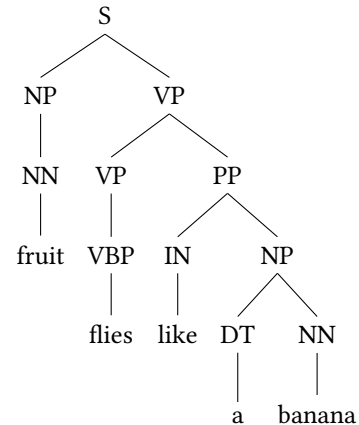


But the grammar also allows other structures, which would lead to other meanings:

(6.51)



(6.52)



Interpretation (6.51) says that a certain kind of fly, the time fly, is fond of arrows. Interpretation (6.52) says that fruits generally fly in the same way that bananas fly.

**Question 4.** The English word *buffalo* has two meanings: it can be a noun (the name of several species of oxen) or a verb (to overpower, overawe, or constrain by superior force or influence; to outwit, perplex). Also, the plural of the noun *buffalo* is *buffalo*. Therefore, the following strings are all grammatical:

(6.53) Buffalo! (Overpower!)

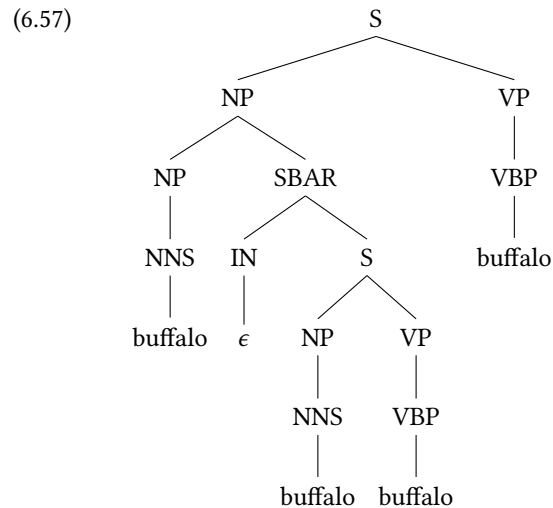
(6.54) Buffalo buffalo. (Oxen overpower.)

(6.55) Buffalo buffalo buffalo. (Oxen overpower oxen.)



(6.56) Buffalo buffalo buffalo buffalo. (Oxen that overpower oxen, overpower.)

In fact, the entire set  $\{\text{buffalo}^n \mid n \geq 1\}$  is a subset of English. Can you write a CFG that generates it according to English grammar? Hint: here is a tree for the fourth example:



How many structures can you find for the following sentence:

(6.58) Buffalo buffalo buffalo buffalo buffalo.

## 6.2.4 Weighted context free grammars

*Weighted CFGs* are a straightforward extension of CFGs. Recall that FSTs map an input string to a set of possible output strings, whereas weighted FSTs give us a *distribution* over possible output strings. In the same way, weighted CFGs help us deal with ambiguity (a single string having multiple structures) by giving us a *distribution* over possible structures.

In a weighted CFG, every production has a weight attached to it, which we write as

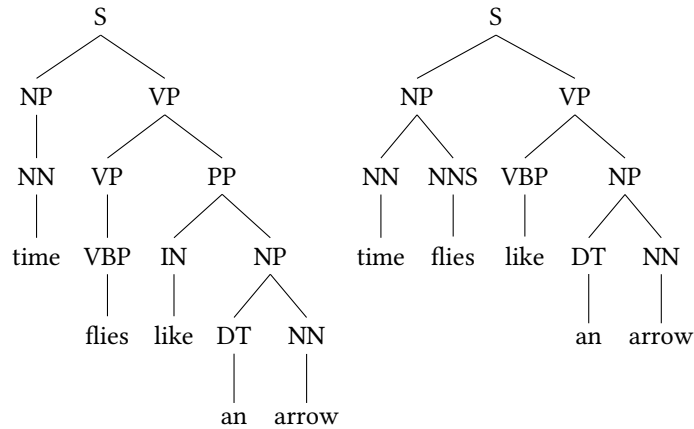
$$A \xrightarrow{p} \beta$$

The weight of a derivation is the product of the weights of the rules used in the derivation (if a rule is used  $k$  times, we multiply its weight in  $k$  times).

Thus we can take our grammar from last time and add weights:

$$\begin{array}{ll}
 S \xrightarrow{1} NP VP & DT \xrightarrow{0.5} a \\
 NP \xrightarrow{0.5} DT NN & DT \xrightarrow{0.5} an \\
 NP \xrightarrow{0.4} NN & NN \xrightarrow{0.25} time \\
 NP \xrightarrow{0.1} NN NNS & NN \xrightarrow{0.25} fruit \\
 VP \xrightarrow{0.6} VBP NP & NN \xrightarrow{0.25} arrow \\
 VP \xrightarrow{0.3} VBP & NN \xrightarrow{0.25} banana \\
 VP \xrightarrow{0.1} VP PP & NNS \xrightarrow{1} flies \\
 PP \xrightarrow{1} IN NP & VBP \xrightarrow{0.5} flies \\
 & VBP \xrightarrow{0.5} like \\
 & IN \xrightarrow{1} like
 \end{array} \tag{6.59}$$

**Question.** What would the weight of these two derivations be?



A *probabilistic CFG* or *PCFG* is one in which the probabilities of all rules with a given left-hand side sum to one (Booth and Thompson, 1973). A PCFG is called *consistent* if the probabilities of all derivations sum to one.

Aren't all PCFGs consistent? Actually, no:

$$S \xrightarrow{0.9} SS \tag{6.60}$$

$$S \xrightarrow{0.1} a \tag{6.61}$$

Let  $P_n$  be the total weight of trees of height  $\leq n$ . Thus

$$P_1 = 0.1 \tag{6.62}$$

$$P_{n+1} = 0.9P_n^2 + 0.1 \tag{6.63}$$

The second equation is because a tree of height  $\leq n + 1$  can either be a tree of height 1, formed using rule (6.61), or a tree formed using rule (6.60) and two trees of height  $\leq n$ . The limit  $P = \lim_{n \rightarrow \infty} P_n$  must be a fixed point of the second equation. There are two fixed points:  $P = 0.9P^2 + 0.1 \Rightarrow P = \frac{1}{9}$  or 1. But note that  $P_1 < \frac{1}{9}$ , and  $P_i < \frac{1}{9} \Rightarrow P_{i+1} < \frac{1}{9}$ . Since the sequence is always less than  $\frac{1}{9}$ , it cannot converge to 1. Therefore,  $P = \frac{1}{9}$ !

**Question.** What happened to the other  $\frac{8}{9}$ ?

## 6.3 Parsing Algorithms

Next, we explore the *parsing* problem, which encompasses several questions, including:

- Does  $L(G)$  contain  $w$ ?
- What is the highest-weight derivation of  $w$ ?
- What is the set of all derivations of  $w$ ?

### 6.3.1 Chomsky normal form

Let's assume that  $G$  has a particularly simple form. We say that a CFG is in *Chomsky normal form* if each of its productions has one of the following forms:

$$\begin{aligned} X &\rightarrow YZ \\ X &\rightarrow a \end{aligned}$$

It can be shown (see below) that any context-free grammar not generating a language containing  $\epsilon$  can be converted into Chomsky normal form, and still generate the same language.

Our grammar from above can be massaged to be in Chomsky normal form:

$$\begin{aligned} S &\rightarrow NP VP \\ NP &\rightarrow DT NN \\ NP &\rightarrow \text{time} \mid \text{fruit} \\ NP &\rightarrow NN NNS \\ VP &\rightarrow VBP NP \\ VP &\rightarrow \text{flies} \\ VP &\rightarrow VP PP \\ PP &\rightarrow IN NP \\ DT &\rightarrow a \mid \text{an} \\ NN &\rightarrow \text{time} \mid \text{fruit} \mid \text{arrow} \mid \text{banana} \\ NNS &\rightarrow \text{flies} \\ VBP &\rightarrow \text{like} \\ IN &\rightarrow \text{like} \end{aligned} \tag{6.64}$$

### 6.3.2 The CKY algorithm

The *CKY algorithm* is named after three people who independently invented it: Cocke, Kasami, and Younger, although it has been rediscovered more times than that.

In its most basic form, the algorithm just decides whether  $w \in L(G)$ . It builds a data structure known as a *chart*; it is an  $n \times n$  array. The element  $chart[i, j]$  is a set of nonterminal symbols. If  $X \in chart[i, j]$ , then that means we have discovered that  $X \Rightarrow^* w_{i+1} \cdots w_j$ .

**Require:** string  $w = w_1 \cdots w_n$  and grammar  $G = (N, \Sigma, R, S)$

**Ensure:**  $w \in L(G)$  iff  $S \in chart[0, n]$

```

1: initialize  $chart[i, j] \leftarrow \emptyset$  for all  $0 \leq i < j \leq n$ 
2: for all  $i \leftarrow 1, \dots, n$  and  $(X \rightarrow w_i) \in R$  do
3:    $chart[i - 1, i] \leftarrow chart[i - 1, i] \cup \{X\}$ 
4: end for
5: for  $\ell \leftarrow 2, \dots, n$  do
6:   for  $i \leftarrow 0, \dots, n - \ell$  do
7:      $j \leftarrow i + \ell$ 
8:     for  $k \leftarrow i + 1, \dots, j - 1$  do
9:       for all  $(X \rightarrow YZ) \in R$  do
10:        if  $Y \in chart[i, k]$  and  $Z \in chart[k, j]$  then
11:           $chart[i, j] \leftarrow chart[i, j] \cup \{X\}$ 
12:        end if
13:      end for
14:    end for
15:  end for
16: end for

```

**Question 5.** What is the time and space complexity of this algorithm?

**Question 6.** Using the grammar (6.64), run the CKY algorithm on the string:

<sub>0</sub> time <sub>1</sub> flies <sub>2</sub> like <sub>3</sub> an <sub>4</sub> arrow <sub>5</sub>

0,1	0,2	0,3	0,4	0,5
	1,2	1,3	1,4	1,5
		2,3	2,4	2,5
			3,4	3,5
				4,5

### 6.3.3 Viterbi CKY

But it is much more useful to find the highest-weight parse. Suppose that our grammar has the following probabilities:

$$\begin{array}{ll}
 S \xrightarrow{1} NP VP & DT \xrightarrow{0.5} a \\
 NP \xrightarrow{0.5} DT NN & DT \xrightarrow{0.5} an \\
 NP \xrightarrow{0.2} time & NN \xrightarrow{0.25} time \\
 NP \xrightarrow{0.2} fruit & NN \xrightarrow{0.25} fruit \\
 NP \xrightarrow{0.1} NN NNS & NN \xrightarrow{0.25} arrow \\
 VP \xrightarrow{0.6} VBP NP & NN \xrightarrow{0.25} banana \\
 VP \xrightarrow{0.3} flies & NNS \xrightarrow{1} flies \\
 VP \xrightarrow{0.1} VP PP & VBP \xrightarrow{1} like \\
 PP \xrightarrow{1} IN NP & IN \xrightarrow{1} like
 \end{array} \tag{6.65}$$

Then we use a modification of CKY that is analogous to the Viterbi algorithm. First, we modify the algorithm to find the maximum weight:

**Require:** string  $w = w_1 \cdots w_n$  and grammar  $G = (N, \Sigma, R, S)$   
**Ensure:**  $best[0, n][S]$  is the maximum weight of a parse of  $w$

- 1: initialize  $best[i, j][X] \leftarrow 0$  **for all**  $0 \leq i < j \leq n, X \in N$
- 2: **for all**  $i \leftarrow 1, \dots, n$  **and**  $(X \xrightarrow{p} w_i) \in R$  **do**
- 3:      $best[i - 1, i][X] \leftarrow \max\{best[i - 1, i][X], p\}$
- 4: **end for**
- 5: **for**  $\ell \leftarrow 2, \dots, n$  **do**
- 6:     **for**  $i \leftarrow 0, \dots, n - \ell$  **do**
- 7:          $j \leftarrow i + \ell$
- 8:         **for**  $k \leftarrow i + 1, \dots, j - 1$  **do**
- 9:             **for all**  $(X \xrightarrow{p} YZ) \in R$  **do**
- 10:                  $p' \leftarrow p \times best[i, k][Y] \times best[k, j][Z]$
- 11:                  $best[i, j][X] \leftarrow \max\{best[i, j][X], p'\}$
- 12:             **end for**
- 13:         **end for**
- 14:     **end for**
- 15: **end for**

**Question 7.** Do you see how to modify the algorithm to compute the *total* weight of all parses of  $w$ ?

A slight further modification lets us find the maximum-weight parse itself. Just as in the Viterbi algorithm for FSAs, whenever we update  $best[i, j][X]$  to a new best weight, we also need to store a *back-pointer* that records how we obtained that weight. We will represent back-pointers like this:  $X_{i,j} \rightarrow Y_{i,k}Z_{k,j}$  means that we built an  $X$  spanning  $i, j$  from a  $Y$  spanning  $i, k$  and a  $Z$  spanning  $k, j$ .

**Require:** string  $w = w_1 \cdots w_n$  and grammar  $G = (N, \Sigma, R, S)$   
**Ensure:**  $G'$  generates the best parse of  $w$   
**Ensure:**  $best[0, n][S]$  is its weight

- 1: **for all**  $0 \leq i < j \leq n, X \in N$  **do**
- 2:     initialize  $best[i, j][X] \leftarrow 0$
- 3:     initialize  $back[i, j][X] \leftarrow \text{nil}$
- 4: **end for**
- 5: **for all**  $i \leftarrow 1, \dots, n$  **and**  $(X \xrightarrow{p} w_i) \in R$  **do**
- 6:     **if**  $p > best[i - 1, i][X]$  **then**
- 7:          $best[i - 1, i][X] \leftarrow p$
- 8:          $back[i - 1, i][X] \leftarrow (X_{i-1,i} \rightarrow w_i)$
- 9:     **end if**
- 10: **end for**
- 11: **for**  $\ell \leftarrow 2, \dots, n$  **do**
- 12:     **for**  $i \leftarrow 0, \dots, n - \ell$  **do**
- 13:          $j \leftarrow i + \ell$
- 14:         **for**  $k \leftarrow i + 1, \dots, j - 1$  **do**
- 15:             **for all**  $(X \xrightarrow{p} YZ) \in R$  **do**
- 16:                  $p' \leftarrow p \times best[i, k][Y] \times best[k, j][Z]$
- 17:                 **if**  $p' > best[i, j][X]$  **then**

```

18:             best[i, j][X] ← p'
19:             back[i, j][X] ← (Xi,j → Yi,kZk,j)
20:         end if
21:     end for
22: end for
23: end for
24: end for
25: G' = {back[i, j][X] | 0 ≤ i < j ≤ n, X ∈ N}

```

$G'$  is then a grammar that generates at most one tree, the best tree for  $w$ .

**Question 8.** Using the grammar (6.65), run the Viterbi CKY algorithm on the same string:

0 time 1 flies 2 like 3 an 4 arrow 5

0,1	0,2	0,3	0,4	0,5
	1,2	1,3	1,4	1,5
		2,3	2,4	2,5
			3,4	3,5
				4,5

**Partial solution.** Cell (0, 5) is missing.

0,1 NP/0.2/NP → time NN/0.25/NN → time	0,2 S/0.06/S → NP <sub>1</sub> VP NP/0.025/NP → NN <sub>1</sub> NNS	0,3 ∅	0,4 ∅	0,5
	1,2 VP/0.3/VP → flies NNS/1/NNS → flies	1,3 ∅	1,4 ∅	1,5 VP/0.001875/VP → VP <sub>2</sub> PP
		2,3 VBP/1/VBP → like IN/1/IN → like	2,4 ∅	2,5 VP/0.0375/VP → VBP <sub>3</sub> NP PP/0.0625/PP → IN <sub>3</sub> NP
			3,4 DT/0.5/IN → an	3,5 NP/0.0625/NP → DT <sub>4</sub> NN
				4,5 NN/0.25/NN → arrow

### 6.3.4 Parsing general CFGs

Previously, we learned about PCFGs, and how to find the best PCFG derivation of a string using the Viterbi algorithm. Now we will extend those algorithms to the general CFG case.

#### Binarization

It turns out that any CFG (whose language does not contain  $\epsilon$ ) can be converted into an equivalent grammar in Chomsky normal form.

To guarantee that  $k \leq 2$ , we must eliminate all rules with right-hand side longer than 2. We will see below that the grammars we extract from training data may already have this property. But if not, we need to *binarize* the grammar. For example, suppose we have the production

$$\text{NP} \rightarrow \text{DT JJS NN NN PP} \quad (6.66)$$

which is too long to be in Chomsky normal form. There are many ways to break this down into smaller rules, but here is one way. We create a bunch of new



nonterminal symbols  $NP(\beta)$  where  $\beta$  is a string of nonterminal symbols; this stands for a partial NP whose sisters to the *left* are  $\beta$ . Then we replace rule (6.66) with:

$$NP \rightarrow DT NP(DT) \quad (6.67)$$

$$NP(DT) \rightarrow JJS NP(DT,JJS) \quad (6.68)$$

$$NP(DT,JJS) \rightarrow NN NP(DT,JJS,NN) \quad (6.69)$$

$$NP(DT,JJS,NN) \rightarrow NN NP(DT,JJS,NN,NN) \quad (6.70)$$

$$NP(DT,JJS,NN,NN) \rightarrow PP \quad (6.71)$$

Note that the annotations contain enough information to reverse the binarization. So the binarized grammar is equivalent to the unbinarized grammar, but has  $k \leq 2$ .

### Parsing with unary rules

But we are not done yet. CKY does not just require  $k \leq 2$ , but also forbids rules of any of the following forms:

$$A \rightarrow ab \quad (6.72)$$

$$A \rightarrow aB \quad (6.73)$$

$$A \rightarrow Ab \quad (6.74)$$

$$A \rightarrow \epsilon \quad (6.75)$$

$$A \rightarrow B \quad (6.76)$$

The first three cases are very easy to eliminate, but we never see them in grammars induced from the Penn Treebank. *Nullary rules* (6.75) are not hard to eliminate (Hopcroft and Ullman, 1979), but the weighted case can be nasty (Stolcke, 1995). Fortunately, nullary rules aren't very common in practice, so we won't bother with them here.

*Unary rules* (6.76) are quite common and annoying. Like nullary rules, they are not hard to eliminate from a CFG (Hopcroft and Ullman, 1979), but in practice, most people don't try to; instead, they extend the CKY algorithm to handle them directly. The extension shown below is not the most efficient, but fits most naturally with the way we have implemented CKY.

**Require:** string  $w = w_1 \cdots w_n$  and grammar  $G = (N, \Sigma, R, S)$

**Ensure:**  $w \in L(G)$  iff  $S \in \text{chart}[0, n]$

- 1: initialize  $\text{chart}[i, j] \leftarrow \emptyset$  **for all**  $0 \leq i < j \leq n$
- 2: **for all**  $i \leftarrow 1, \dots, n$  **and**  $(X \rightarrow w_i) \in R$  **do**
- 3:      $\text{chart}[i - 1, i] \leftarrow \text{chart}[i - 1, i] \cup \{X\}$
- 4: **end for**
- 5: **for**  $\ell \leftarrow 2, \dots, n$  **do**
- 6:     **for**  $i \leftarrow 0, \dots, n - \ell$  **do**
- 7:          $j \leftarrow i + \ell$
- 8:         **for**  $k \leftarrow i + 1, \dots, j - 1$  **do**
- 9:             **for all**  $(X \rightarrow YZ) \in R$  **do**
- 10:                 **if**  $Y \in \text{chart}[i, k]$  **and**  $Z \in \text{chart}[k, j]$  **then**

```

11:         chart[i, j] ← chart[i, j] ∪ {X}
12:     end if
13: end for
14: end for
15: again ← true
16: while again do
17:     again ← false
18:     for all (X → Y) ∈ R do
19:         if X ∉ chart[i, j] and Y ∈ chart[i, j] then
20:             chart[i, j] ← chart[i, j] ∪ {X}
21:             again ← true
22:         end if
23:     end for
24: end while
25: end for
26: end for

```

The new part is lines 15–24 and is analogous to the binary rule case.

**Question.** Why is the **while** loop on line 16 necessary? What is its maximum number of iterations?

Finally, let's put together weights (Viterbi CKY) and unary rules:

**Require:** string  $w = w_1 \cdots w_n$  and grammar  $G = (N, \Sigma, R, S)$

**Ensure:**  $G'$  generates the best parse of  $w$

**Ensure:**  $best[0, n][S]$  is its weight

```

1: for all  $0 \leq i < j \leq n, X \in N$  do
2:     initialize  $best[i, j][X] \leftarrow 0$ 
3:     initialize  $back[i, j][X] \leftarrow \text{nil}$ 
4: end for
5: for all  $i \leftarrow 1, \dots, n$  and  $(X \xrightarrow{p} w_i) \in R$  do
6:     if  $p > best[i - 1, i][X]$  then
7:          $best[i - 1, i][X] \leftarrow p$ 
8:          $back[i - 1, i][X] \leftarrow (X_{i-1, i} \rightarrow w_i)$ 
9:     end if
10: end for
11: for  $\ell \leftarrow 2, \dots, n$  do
12:     for  $i \leftarrow 0, \dots, n - \ell$  do
13:          $j \leftarrow i + \ell$ 
14:         for  $k \leftarrow i + 1, \dots, j - 1$  do
15:             for all  $(X \xrightarrow{p} YZ) \in R$  do
16:                  $p' \leftarrow p \times best[i, k][Y] \times best[k, j][Z]$ 
17:                 if  $p' > best[i, j][X]$  then
18:                      $best[i, j][X] \leftarrow p'$ 
19:                      $back[i, j][X] \leftarrow (X_{i, j} \rightarrow Y_{i, k} Z_{k, j})$ 
20:                 end if
21:             end for
22:         end for

```

```

23:   again ← true
24:   while again do
25:     again ← false
26:     for all  $(X \xrightarrow{p} Y) \in R$  do
27:        $p' \leftarrow p \times \text{best}[i, j][Y]$ 
28:       if  $p' > \text{best}[i, j][X]$  then
29:          $\text{best}[i, j][X] = p'$ 
30:          $\text{back}[i, j][X] \leftarrow (X_{i,j} \rightarrow Y_{i,j})$ 
31:         again ← true
32:       end if
33:     end for
34:   end while
35: end for
36: end for
37:  $G' \leftarrow \text{EXTRACT}(S, 0, n)$ 

```

If the grammar has unary cycles in it, that is, it is possible to derive  $X \Rightarrow \dots \Rightarrow^* X$ , then certain complications can arise from the fact that a string may have an infinite number of derivations. In particular, if the weight of the cycle is greater than 1, then the Viterbi CKY algorithm will break. Even if all rule weights are less than 1, some algorithms require modification; for example, if we want to find the total weight of all the derivations of a string, we have to perform an infinite summation (Stolcke, 1995). Therefore, it is fairly common to implement hacks of various kinds to break the cycles. For example, we could modify the grammar so that it goes round the cycle at most five times.

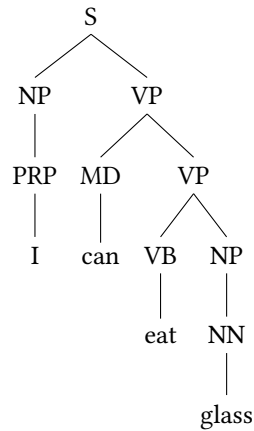
**Question 9.** Why doesn't the Viterbi CKY algorithm break on unary cycles if we assume that all rule weights are less than 1?

## 6.4 Neural Parsing

The literature on parsing can be divided into *dependency parsing* and *constituency parsing*. The former is probably more widespread, but we're focusing on the latter. Neural constituency parsing models can be further divided into models that are (loosely) based on pushdown automata (PDAs) and models that are based on context-free grammars.

### 6.4.1 Stack-based approaches

These approaches treat trees as sequences of symbols and use some kind of stack to ensure that the sequence is well-formed. I give just one example, which is to train a neural language model on trees represented as bracketed strings (Vinyals et al., 2015; Choe and Charniak, 2016). For example,



would become something like

[S [NP [PRP I PRP] NP] [VP [MD can MD] [VP [VB eat VB] [NP [NN glass NN] NP] VP] VP] S]

To parse, we search among flattened trees for the one which has the highest probability according to the language model. This is an intractable search, but we can make it practical using beam search, in which at each time step we keep only the top  $k$  possibilities.

The neural language model is a RNN. I'm not aware of approaches that use Transformer language models, but that would certainly be possible.

## 6.4.2 Grammar-based approaches

Another approach is to use a weighted CFG, but to compute the weights of the CFG using a neural network. Suppose that the probability of a tree is still

$$P(\text{tree}) = \prod_{(A \rightarrow \beta) \in \text{tree}} P(A \rightarrow \beta)$$

but imagine that  $P(A \rightarrow \beta)$  is computed by a neural network, in the hopes that it can learn better probabilities than relative frequency estimation ( $P(A \rightarrow \beta) = c(A \rightarrow \beta)/c(A)$ ). However, relative-frequency estimation is already guaranteed to maximize the likelihood of the training data, so it's not possible to do better! The only way to improve the model is to let the rule scores depend on information outside the rule itself.

For example, maybe a rule probability could depend on the rest of the words in the sentence. This doesn't make sense mathematically, because words depend on the rules, so it's circular for the rules to depend on the words.<sup>2</sup>

The way out is to change from a *generative* model, that is, a model of  $P(\text{tree}, w)$  to a *discriminative* model, that is, a model of  $P(\text{tree} | w)$ . Since the tree now depends on the words, instead of the other way around, any rule can depend on any of the words in the sentence.

<sup>2</sup>Although, admittedly, it worked pretty well when I tried it. Something similar happens in the hybrid HMM-DNN model that is still the state of the art in speech recognition.

The model described here is the one provided in HW3. It's a highly simplified version of several neural models that were state-of-the-art at the time (Durrett and Klein, 2015; Stern, Andreas, and Klein, 2017). We start off with an encoder as in a translation system:

$$\begin{aligned} \mathbf{V} &\in \mathbb{R}^{n \times d} \\ \mathbf{V}_i &= \text{Embedding}^{\boxed{1}}(w_i) & i = 1, \dots, n \end{aligned} \quad (6.77)$$

$$\begin{aligned} \mathbf{H} &\in \mathbb{R}^{n \times d} \\ \mathbf{H} &= \text{RNN}^{\boxed{2}}(\mathbf{V}). \end{aligned} \quad (6.78)$$

The encoder could just as easily be a Transformer encoder (Kitaev and Klein, 2018).

If this were a language model, you'd expect something like a softmax layer to predict the next word. But we want to predict a whole tree, so we have something like a softmax over trees:

$$P(\text{tree} \mid \mathbf{w}) = \frac{\exp s(\text{tree})}{\sum_{\text{tree}'} \exp s(\text{tree}')} \quad (6.79)$$

where the summation over  $\text{tree}'$  sums over all trees with  $\mathbf{w}$  on their leaves, and  $s(\cdot)$  is the score of a tree, which is the sum of the scores of the rules in it:

$$s(\text{tree}) = \sum_{(A \rightarrow \beta) \in \text{tree}} s(A \rightarrow \beta). \quad (6.80)$$

This  $s(A \rightarrow \beta)$  is the score of rule  $A \rightarrow \beta$ , where  $A$  spans the substring  $w_{i+1} \cdots w_j$ . It is computed by a feedforward neural network. Define a new kind of module,

$$\text{LinearLayer}^{\boxed{1}}(\mathbf{x}) = \mathbf{W}^{\boxed{1}} \mathbf{x} + \mathbf{b}^{\boxed{1}}. \quad (6.81)$$

Then for each span  $(i, j)$ , compute a *span encoding*

$$\mathbf{s}^{(i,j)} = \begin{bmatrix} \mathbf{H}_{i+1,:} \\ \mathbf{H}_{j,:} \end{bmatrix} \quad (6.82)$$

and then feed it through two layers:

$$\begin{aligned} \mathbf{h}^{(i,j)} &\in \mathbb{R}^d \\ \mathbf{h}^{(i,j)} &= \tanh(\text{LinearLayer}^{\boxed{3}}(\mathbf{s}^{(i,j)})) \end{aligned} \quad (6.83)$$

$$\begin{aligned} \mathbf{o}^{(i,j)} &\in \mathbb{R}^{|G|} \\ \mathbf{o}^{(i,j)} &= \text{LinearLayer}^{\boxed{4}}(\mathbf{h}^{(i,j)}). \end{aligned} \quad (6.84)$$

If we number the rules of the grammar and rule  $A \rightarrow \beta$  has index  $r$ , then the rule score  $s(A \rightarrow \beta)$  is  $\mathbf{o}_r^{(i,j)}$ .

Two challenges remain. First, how do we use this model to search for the best

tree? This turns out to be a simple modification to the CKY algorithm.

$$\arg \max_{\text{tree}} P(\text{tree} \mid w) = \arg \max_{\text{tree}} \frac{\exp s(\text{tree})}{\sum_{\text{tree}'} \exp s(\text{tree}')} \quad (6.85)$$

$$= \arg \max_{\text{tree}} \exp s(\text{tree}) \quad (6.86)$$

$$= \arg \max_{\text{tree}} \prod_{(A \rightarrow \beta) \in \text{tree}} \exp s(A \rightarrow \beta). \quad (6.87)$$

This is the same problem that the Viterbi CKY algorithm solves, except we just change the weight of a rule from  $P(A \rightarrow \beta)$  to  $\exp s(A \rightarrow \beta)$ . The weights don't sum to one like probabilities do, but nothing about CKY depends on whether the weights sum to one, so it will work fine.

The second challenge is, how do we compute  $\sum_{\text{tree}'} \exp s(\text{tree}')$  in the denominator of (6.79)? This is needed to form the loss function during training. It is the sum of the weights of all trees with  $w$  on their leaves. This, too, is a simple modification to the CKY algorithm. It builds a table  $total[i, j]$ , and whenever two items  $X$  are added to the same cell  $total[i, j]$ , instead of keeping the best one, we add them. The pseudocode is as follows:

**Require:** string  $w = w_1 \cdots w_n$  and grammar  $G = (N, \Sigma, R, S)$

**Ensure:**  $total[0, n][S]$  is the log of the total weight of all parses of  $w$

```

1: initialize  $total[i, j][X] \leftarrow -\infty$  for all  $0 \leq i < j \leq n, X \in N$ 
2: for all  $i \leftarrow 1, \dots, n$  do
3:    $a \leftarrow w_i$ 
4:   if  $(X \rightarrow a) \in R$  then
5:      $total[i - 1, i][X] \leftarrow \text{logaddexp}(total[i - 1, i][X], s(X \rightarrow_{i-1} a_i))$ 
6:   end if
7: end for
8: for  $\ell \leftarrow 2, \dots, n$  do
9:   for  $i \leftarrow 0, \dots, n - \ell$  do
10:     $j \leftarrow i + \ell$ 
11:    for  $k \leftarrow i + 1, \dots, j - 1$  do
12:      if  $(X \rightarrow YZ) \in R$  then
13:         $p' \leftarrow s(X \rightarrow_i YZ_j) + total[i, k][Y] + total[k, j][Z]$ 
14:         $total[i, j][X] \leftarrow \text{logaddexp}(total[i, j][X], p')$ 
15:      end if
16:    end for
17:  end for
18: end for

```

The function  $\text{logaddexp}$  is defined as

$$\text{logaddexp}(x, y) = \log(\exp x + \exp y) \quad (6.88)$$

but is implemented in a way that avoids over/underflow.

# Bibliography

- Booth, Taylor L. and Richard A. Thompson (1973). "Applying Probability Measures to Abstract Languages". In: *IEEE Trans. Computers* C-22.5, pp. 442–450.
- Choe, Do Kook and Eugene Charniak (Nov. 2016). "Parsing as Language Modeling". In: *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*. Austin, Texas: Association for Computational Linguistics, pp. 2331–2336. doi: 10.18653/v1/D16-1257. URL: <https://www.aclweb.org/anthology/D16-1257>.
- Durrett, Greg and Dan Klein (July 2015). "Neural CRF Parsing". In: *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Beijing, China: Association for Computational Linguistics, pp. 302–312. doi: 10.3115/v1/P15-1030. URL: <https://www.aclweb.org/anthology/P15-1030>.
- Hopcroft, John E. and Jeffrey D. Ullman (1979). *Introduction to Automata Theory, Languages, and Computation*. Reading, MA: Addison-Wesley.
- Kitaev, Nikita and Dan Klein (July 2018). "Constituency Parsing with a Self-Attentive Encoder". In: *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Melbourne, Australia: Association for Computational Linguistics, pp. 2676–2686. doi: 10.18653/v1/P18-1249. URL: <https://www.aclweb.org/anthology/P18-1249>.
- Stern, Mitchell, Jacob Andreas, and Dan Klein (July 2017). "A Minimal Span-Based Neural Constituency Parser". In: *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Vancouver, Canada: Association for Computational Linguistics, pp. 818–827. doi: 10.18653/v1/P17-1076. URL: <https://www.aclweb.org/anthology/P17-1076>.
- Stolcke, Andreas (1995). "An Efficient Probabilistic Context-Free Parsing Algorithm that Computes Prefix Probabilities". In: *Computational Linguistics* 21, pp. 165–201.
- Vinyals, Oriol et al. (2015). "Grammar as a Foreign Language". In: *Advances in Neural Information Processing Systems*. Ed. by C. Cortes et al. Vol. 28. Curran Associates, Inc. URL: <https://proceedings.neurips.cc/paper/2015/file/277281aada22045c03945dcb2ca6f2ec-Paper.pdf>.