

Context-Free Grammars

A canonical first example of a **context-free grammar (CFG)** is shown in (1). This grammar generates the stringset $\{a^n b^n | n \geq 1\}$, i.e. the stringset $\{ab, aabb, aaabbb, aaaabbbb, \dots\}$.

- (1) $S \rightarrow aSb$
 $S \rightarrow ab$
- Handwritten notes:* n repetitions of a , n repetitions of b

The fact that the grammar in (1) generates the string 'aaaabbbb' is due to the fact that we can start from 'S' and gradually rewrite nonterminals in accord with the rules, one step at a time, as follows:

- (2) S
 aSb
 $aaSbb$
 $aaaSbbb$
 $aaaabbbb$
- Handwritten notes:* $S \rightarrow aSb$ (rewritten 3 times), $S \rightarrow ab$ (rewritten 1 time)

A sequence of strings related to each other as in (2) is known as a **derivation**.

1 Formal definition of context-free grammars

- (3) A **context-free grammar (CFG)** is a four-tuple (N, Σ, I, R) where:
- N is a finite set of nonterminal symbols; *internally*
 - Σ , the alphabet, is a finite set of terminal symbols (disjoint from N); *terminal*
 - $I \subseteq N$ is the set of initial nonterminal symbols; and *top of tree*
 - $R \subseteq (N \times (N \cup \Sigma)^*)$ is a finite set of rules.

Handwritten notes: N is nonterminal, $(N \cup \Sigma)^*$ is sequence of terminal and nonterminal.

Strictly speaking, (1) is an informal representation of the following mathematical object:

- (4) $G = (\{S\}, \{a,b\}, \{S\}, \{(S, aSb), (S, ab)\})$

We will write $u \xrightarrow{G} v$ to mean that we can derive the string v from the string u by rewriting one symbol in u in accordance with the rules of the grammar G . We write $u \xRightarrow{G} v$ to mean that we can derive v from u in zero or more such steps.

So, what we did above in (2) can be described in symbols like this:

- (5) a. $S \xrightarrow{G} aSb \xrightarrow{G} aaSbb \xrightarrow{G} aaaSbbb \xrightarrow{G} aaaabbbb$ *one step*
 b. $S \xRightarrow{G} * aaaabbbb$ *zero or more steps (four)*

A CFG $G = (N, \Sigma, I, R)$ generates a string $w \in \Sigma^*$ iff there is some way to derive w from some initial nonterminal symbol:

$$(6) \quad w \in \mathcal{L}(G) \iff \bigvee_{n \in N} [I(n) \wedge n \xrightarrow{G}^* w]$$

Here are two handy conventions that we'll adopt:

- Unless otherwise specified, we'll generally use uppercase letters for nonterminal symbols and lowercase letters for terminal symbols, which means that we can recover the sets N and Σ if we're only given a collection of rules.
- Unless otherwise specified, we'll assume that there is exactly one initial nonterminal symbol, and that it is the symbol on the left hand side of first rule listed. This way the grammar can be entirely identified by listing its rules.

2 Equivalent, leftmost, and rightmost derivations

Now let's consider a more interesting grammar:

$$(7) \quad \begin{aligned} N &= \{\text{NP}\}, I = \{\text{NP}\} \\ \Sigma &= \{\text{and, or, apples, bananas, oranges}\} \end{aligned}$$

NP \rightarrow NP and NP
 NP \rightarrow NP or NP
 NP \rightarrow apples
 NP \rightarrow bananas
 NP \rightarrow oranges

Notice that the string 'apples and oranges' has two distinct derivations in this grammar:

$$(8) \quad \begin{array}{ll} \text{a.} & \text{NP} \\ & \text{NP and NP} \\ & \text{apples and NP} \\ & \text{apples and oranges} \\ \text{b.} & \text{NP} \\ & \text{NP and NP} \\ & \text{NP and oranges} \\ & \text{apples and oranges} \end{array}$$

And the string 'apples and oranges or bananas' has many distinct derivations, but here are two of them:

$$(9) \quad \begin{array}{ll} \text{a.} & \text{NP} \\ & \text{NP and NP} \\ & \text{apples and NP} \\ & \text{apples and NP or NP} \\ & \text{apples and oranges or NP} \\ & \text{apples and oranges or bananas} \\ \text{b.} & \text{NP} \\ & \text{NP or NP} \\ & \text{NP and NP or NP} \\ & \text{apples and NP or NP} \\ & \text{apples and oranges or NP} \\ & \text{apples and oranges or bananas} \end{array}$$

But in an important way, the two derivations in (8) seem to be just different ways of "doing the same thing," whereas the two in (9) seem substantively different.

The underlying distinction is that the two derivations in (8) are both plugging together the following three facts:

- (10) NP \Rightarrow^* apples and oranges constituent
 NP \Rightarrow^* apples
 NP \Rightarrow^* oranges

But the two derivations in (9) plug together different collections of facts, of this sort:

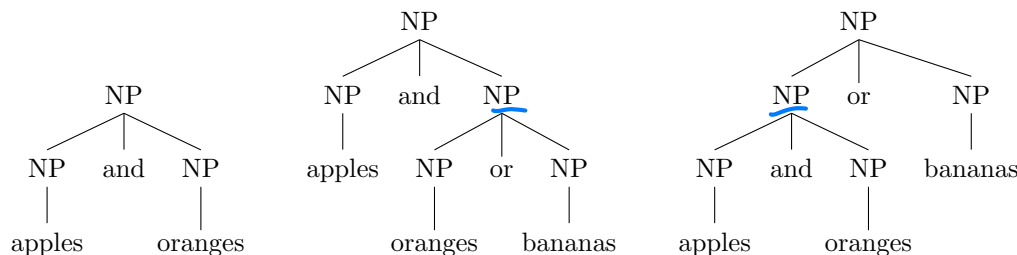
- (11) a. NP \Rightarrow^* apples and oranges or bananas
 NP \Rightarrow^* apples
 NP \Rightarrow^* oranges or bananas
 NP \Rightarrow^* oranges
 NP \Rightarrow^* bananas
- b. NP \Rightarrow^* apples and oranges or bananas
 NP \Rightarrow^* apples and oranges
 NP \Rightarrow^* apples
 NP \Rightarrow^* oranges
 NP \Rightarrow^* bananas

In the early generative linguistics literature, these individual constituency-facts like ‘NP \Rightarrow^* oranges’ or ‘NP \Rightarrow^* bananas’ were known as *is-a* relationships (e.g. ‘oranges or bananas’ *is a* NP).¹

Let’s call such a collection of *is-a* relationships an **analysis** of a string. So:

- (8)a and (8)b both correspond to the *same analysis of the string* ‘apples and oranges,’ shown in (10).
- (9)a and (9)b correspond to two *distinct analyses of the string* ‘apples and oranges or bananas,’ shown in (11)a and (11)b, respectively.

An analysis can be represented graphically by a tree:



If two derivations correspond to the *same analysis*, i.e. they both express the same set of *is-a* relations, we’ll call them **equivalent derivations**.

The relationship between derivations and analyses is many-to-one: usually we only care about distinct *analyses* of a string, and the differences between equivalent derivations like (8)a and (8)b concern only the order in which the rewrites take place. We can get rid of these “spurious choices” by adopting some fixed strategy that dictates, given any intermediate point in a derivation (such as ‘NP and NP’), which nonterminal symbol must be rewritten next.

- A **leftmost derivation** is a derivation where every step rewrites the leftmost nonterminal symbol in the preceding string. There is exactly one leftmost derivation per analysis.
- A **rightmost derivation** is a derivation where every step rewrites the rightmost nonterminal symbol in the preceding string. There is exactly one rightmost derivation per analysis.

The derivation in (8)a is a leftmost derivation; the derivation in (8)b is a rightmost derivation. The two derivations in (9) are both leftmost derivations.

¹See e.g. chapter 4 of Chomsky’s *Syntactic Structures* (1957), or chapter 1 of Lasnik’s *Syntactic Structures Revisited* (2000).

By considering only leftmost derivations or considering only rightmost derivations, we can leave aside the “irrelevant” differences; if we find two leftmost derivations or two rightmost derivations for a string, then it has two distinct tree structures, i.e. can be analyzed by two distinct sets of *is-a* relationships.

3 Finding analyses

We'll assume from here on that all CFGs are in **Chomsky Normal Form (CNF)**. In CNF, every rule has on its right-hand side either:

- a single terminal symbol, or
- exactly two nonterminal symbols.

So in effect, we'll take R to be of the form $R \subseteq N \times ((N \times N) \cup \Sigma)$.

Here's an example grammar in CNF:

(12) $N = \{VP, NP, PP, V, P\}$
 $\Sigma = \{\text{watches, spies, telescopes, with}\}$
 $I = \{VP\}$

VP \rightarrow V NP
 VP \rightarrow VP PP
 NP \rightarrow NP PP
 PP \rightarrow P NP
 VP \rightarrow spies
 VP \rightarrow watches

V \rightarrow watches
 NP \rightarrow watches
 NP \rightarrow spies
 NP \rightarrow telescopes
 P \rightarrow with

Handwritten notes:
 - N : nonterminal
 - Σ : alphabet
 - I : initial
 - VP : nonterminal
 - $VP \rightarrow V NP$: or
 - $VP \rightarrow VP PP$: or
 - $NP \rightarrow NP PP$: or
 - $PP \rightarrow P NP$: or
 - $VP \rightarrow spies$: terminal
 - $VP \rightarrow watches$: terminal
 - $V \rightarrow watches$: terminal
 - $NP \rightarrow watches$: terminal
 - $NP \rightarrow spies$: terminal
 - $NP \rightarrow telescopes$: terminal
 - $P \rightarrow with$: terminal

Assuming this form is computationally convenient and doesn't really restrict what we can do. Any CFG can be converted to a CFG in CNF that generates the same string set, modulo the empty string.²

The important idea about how a CFG generates a string can be stated in the same form as we saw for FSAs:

$$(13) \quad w \in \mathcal{L}(G) \Leftrightarrow \bigvee_{\text{all possible trees } t} [\text{string } w \text{ can be generated via tree } t]$$

For an FSA, the relevant notion of a path is a linear path through the states; here, think of a tree as a branching path through the nonterminals.

- One way in which CFGs are more complicated than FSAs is that, even given a particular string we're interested in, we don't know what the shape(s) of the relevant path(s) might be.
- So, spelling out what it means to consider “all possible trees t ” is a bit trickier than spelling out “all possible paths p ” for an FSA.

But if we leave aside the “middle” of the tree for the moment, and just focus on the top and the bottom, we can at least write the following:

$$(14) \quad x_1 x_2 \dots x_n \in \mathcal{L}(G) \Leftrightarrow \bigvee_{r \in N} \dots \bigvee_{c_1 \in N} \dots \bigvee_{c_n \in N} [I(r) \wedge \dots \wedge R(c_1, x_1) \wedge R(c_2, x_2) \dots \wedge R(c_n, x_n)]$$

Handwritten notes:
 - $r \in N$: across all nonterminal symbols
 - c_1, c_2, \dots, c_n : and c_1, c_2, \dots, c_n

²That is, for any CFG G we can construct another CFG G' such that G' is in CNF and $\mathcal{L}(G') = \mathcal{L}(G) - \{\epsilon\}$.

For example, what we would need to evaluate to see whether the grammar in (12) will generate ‘watches with telescopes’ would look something like this (still leaving aside the “middle” of the tree):

$$(15) \quad \bigvee_{r \in N} \cdots \bigvee_{c_1 \in N} \bigvee_{c_2 \in N} \bigvee_{c_3 \in N} \left[I(r) \wedge \cdots \wedge R(c_1, \text{watches}) \wedge R(c_2, \text{with}) \wedge R(c_3, \text{telescopes}) \right]$$

terminating symbols

Interleaving the growing of structure with the checking of this structure against the grammar, analogously to what we did with forward and backward values for FSAs, will let us sidestep the issue of working out exactly what would need to go into this “global disjunction.”

4 Inside values

For any CFG G there's a two-place predicate $inside_G$, relating strings to nonterminals.

we can rewrite through series of steps"

string nonterminal

$$(16) \quad inside_G(w)(n) \text{ is true iff the string } w \text{ is derivable from the nonterminal } n \text{ in grammar } G, \text{ i.e. } n \xRightarrow{G}^* w$$

Given a way to work out $inside_G(w)(n)$ for any string and any nonterminal, we can easily use this to check for membership in $\mathcal{L}(G)$:

start and inside
:=
true

$$(17) \quad w \in \mathcal{L}(G) \Leftrightarrow \bigvee_{n \in N} \left[I(n) \wedge inside_G(w)(n) \right]$$

w = x₁ ... x_n

For all nonterminals n

we can do

n

w = x₁ ... x_n

We can arrange the values of the predicate $inside_G$ in a table or chart. The chart in (18) shows the values of $inside_G$ for all non-empty substrings, or *infixes*, of ‘watches spies with telescopes’, where G is the grammar in (12). Each cell in the chart corresponds to one of these infixes.³

(18)

inside value of watch-s

... watches ... spies ... with ... telescopes

VP: 1	VP: 1	VP: 0	VP: 1
NP: 1	NP: 0	NP: 0	NP: 0
PP: 0	PP: 0	PP: 0	PP: 0
V: 1	V: 0	V: 0	V: 0
P: 0	P: 0	P: 0	P: 0
VP: 1	VP: 0	VP: 1	
NP: 1	NP: 0	NP: 1	
PP: 0	PP: 0	PP: 0	
V: 0	V: 0	V: 0	
P: 0	P: 0	P: 0	
VP: 0	VP: 0	VP: 0	
NP: 0	NP: 0	NP: 0	
PP: 0	PP: 1		
V: 0	V: 0		
P: 1	P: 0		
VP: 0			
NP: 1			
PP: 0			
V: 0			
P: 0			

Is there up → watches

is there NP → watches

VP → spies

NP → spies

spies with telescopes

then

start

split up string

$x_1 || \dots || x_m$

x_1 $x_2 x_3 \dots x_m$

We can fill this table in gradually by “working from small to big,” similarly to what we saw for forward and backward values for FSAs, but with some important differences.

- Here it’s the cells on the diagonal, corresponding to infixes of length one, that can be filled in just by lookups into the grammar.⁴
- We then work upwards and to the right to fill in values for longer infixes, which can be calculated by looking just “one step back” at values for shorter infixes.
- But exactly what it means to look “one step back” is a bit more complicated because of the fact that we’re dealing with all infixes, not just prefixes and suffixes.

Specifically, every cell in the table can be filled in according to the following recursive definition (which is not as complicated as it looks):

(19)

a. $inside_G(x)(n) = R(n, x)$

b. $inside_G(x_1 \dots x_m)(n) = \bigvee_{1 \leq i < m} \bigvee_{\ell \in N} \bigvee_{r \in N} [R(n, \ell, r) \wedge inside_G(x_1 \dots x_i)(\ell) \wedge inside_G(x_{i+1} \dots x_m)(r)]$

Does rule exist?

half from nonterminal

This algorithm for filling in a chart with inside values is known as the **CKY algorithm** (after Cocke, Kasami

³Since each cell in this chart corresponds to an infix (and specifies a value for each nonterminal), each cell is analogous to a *column* in the tables we saw earlier for forward and backward values.

⁴Like the leftmost column, representing the empty string, for forward values, and the rightmost column, representing the empty string, for backward values.

and Younger, who all invented/discovered it independently in the late 1960s).

It's interesting to compare the definition in (19) with the definition we saw for backward values.

$$(20) \quad \begin{aligned} \text{a. } bwd_M(\epsilon)(q) &= F(q) \\ \text{b. } bwd_M(x_1 \dots x_n)(q) &= \bigvee_{q_1 \in Q} [\Delta(q, x_1, q_1) \wedge bwd_M(x_2 \dots x_n)(q_1)] \end{aligned}$$

The difference is that:

- a backward value for $x_1 \dots x_n$ is broken down in terms of a single backward value, for $x_2 \dots x_n$ (plus a transition step), whereas
- an inside value for $x_1 \dots x_m$ is broken down in terms of two other inside values, for $x_1 \dots x_i$ and for $x_{i+1} \dots x_m$ (plus a transition step).

So only suffixes of a string have backward values, whereas any sub-parts of a string can have inside values.

5 Outside values

For any CFG G there's a predicate $outside_G$, which relates *pairs of strings* to nonterminals:

$$(21) \quad outside_G(u, v)(n) \text{ is true iff we can get from an initial symbol to the nonterminal } n \text{ in a way that puts the string } u \text{ on its left and puts the string } v \text{ on its right}$$

Given a way to work out $outside_G(u, v)(n)$ for any pair of strings and any nonterminal, we can easily use this to check for membership in $\mathcal{L}(G)$ — in fact we can do this in many ways, one for each symbol in the string of interest.

$$(22) \quad \text{for any } i \in \{1, \dots, m\} : x_1 \dots x_m \in \mathcal{L}(G) \Leftrightarrow \bigvee_{n \in N} [outside_G(x_1 \dots x_{i-1}, x_{i+1} \dots x_m)(n) \wedge R(n, x_i)]$$

Things are getting a bit more complicated here. An outside value concerns a *pair of strings*, because that's what you get if you take an infix — the kind of thing that has an inside value — away from a string.

- Compared to what we saw with FSAs, this contrasts with the fact that, if you take a prefix or a suffix away from a string, you just get another string.

And because of the asymmetry between inside and outside values, it turns out that computing outside values requires first computing inside values.

- This means that (22) isn't itself a useful way to check well-formedness of a string; one would just use (17) instead.
- But being able to compute outside values is necessary for some common tasks in NLP⁵ and helps us get a deeper understanding of how CFGs actually work.

⁵For example, the “inside-outside” algorithm for estimating the rule probabilities in a probabilistic CFG based on unparsed text. See e.g. Manning & Schütze's *Foundations of Statistical NLP*, Ch. 11.

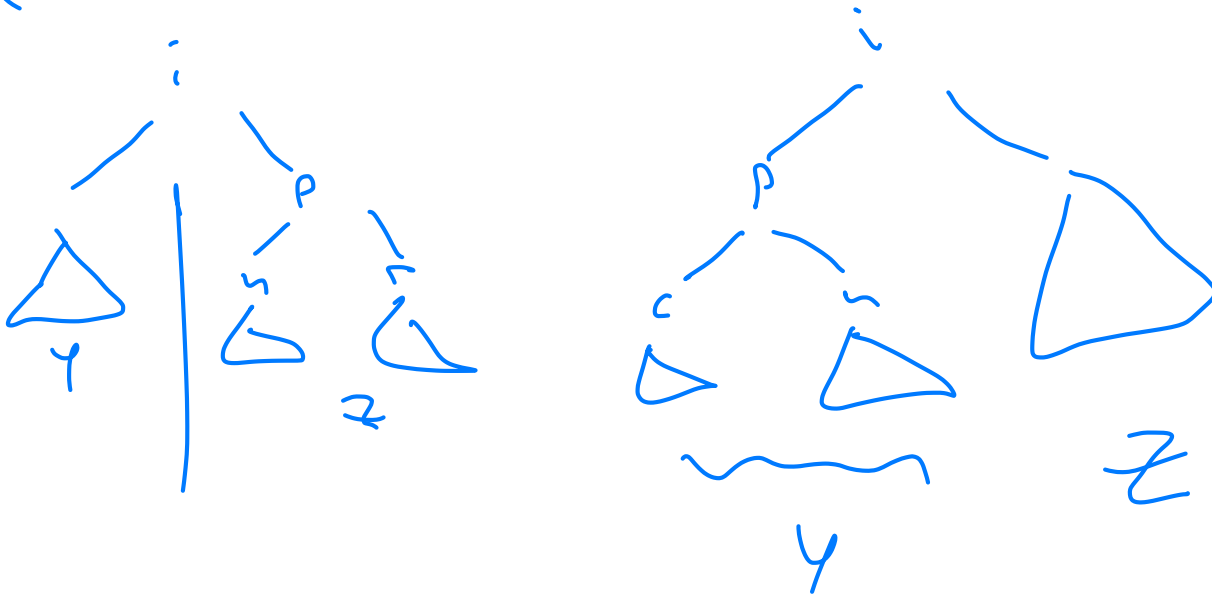
Here's the recursive definition of $outside_G$. It is not as scary as it looks— a picture helps a lot.

- (23) a. $outside_G(\epsilon, \epsilon)(n) = I(n)$
 b. $outside_G(y_1 \dots y_m, z_1 \dots z_q)(n) =$

$$\bigvee_{0 \leq i \leq q} \bigvee_{p \in N} \bigvee_{r \in N} \left[outside_G(y_1 \dots y_m, z_{i+1} \dots z_q)(p) \wedge R(p, n, r) \wedge inside_G(z_1 \dots z_i)(r) \right] \\ \vee \bigvee_{0 \leq i < m} \bigvee_{p \in N} \bigvee_{\ell \in N} \left[outside_G(y_1 \dots y_i, z_1 \dots z_q)(p) \wedge R(p, \ell, n) \wedge inside_G(y_{i+1} \dots y_m)(\ell) \right]$$

The key to understanding (23) is to realize that we need to deal separately with two distinct ways in which we might arrive at the nonterminal symbol n as we work downwards from an initial symbol. Either:

- n was produced as the left daughter of some other parent nonterminal symbol p , via $R(p, n, r)$ for some other right daughter r ; or
- n was produced as the right daughter of some other parent nonterminal symbol p , via $R(p, \ell, n)$ for some other left daughter ℓ .



6 Inside and outside values together

Inside values and outside values can be “snapped together” to describe the generation of a complete string, just like forward and backward values for FSAs:

$$(24) \quad \underbrace{uvw}_{\text{generated by grammar}} \in \mathcal{L}(G) \iff \bigvee_{n \in N} \left[\text{outside}_G(u, w)(n) \wedge \text{inside}_G(v)(n) \right]$$

Some interesting questions:

- Why does this implication only go in one direction? *uvw arbitrary split*
- Why do the implications in the earlier special cases in (17) and (22) go both ways?

CFGs

$$\begin{aligned} uvw \in \mathcal{L}(G) &\iff \bigvee_{n \in N} \left[\text{outside}_G(u, w)(n) \wedge \text{inside}_G(v)(n) \right] \\ w \in \mathcal{L}(G) &\iff \bigvee_{n \in N} \left[I(n) \wedge \text{inside}_G(w)(n) \right] \\ u x w \in \mathcal{L}(G) &\iff \bigvee_{n \in N} \left[\text{outside}_G(u, w)(n) \wedge R(n, x) \right] \end{aligned}$$

FSAs

$$\begin{aligned} uv \in \mathcal{L}(M) &\iff \bigvee_{q \in Q} \left[\text{fwd}_M(u)(q) \wedge \text{bwd}_M(v)(q) \right] \quad \text{all linear} \\ w \in \mathcal{L}(M) &\iff \bigvee_{q \in Q} \left[I(w)(q) \wedge \text{bwd}_M(q) \right] \\ w \in \mathcal{L}(M) &\iff \bigvee_{q \in Q} \left[\text{fwd}_M(w)(q) \wedge F(q) \right] \end{aligned}$$

Also notice how:

- The definition of *outside* makes use of I and $R(c, c_1, c_2)$ rules, but not $R(c, x)$ rules; the definition of *fwd* makes use of I and Δ , but not F . *like forward*
- The definition of *inside* makes use of $R(c, c_1, c_2)$ rules and $R(c, x)$ rules, but not I ; the definition of *bwd* makes use of Δ and F , but not I . *like backwards*