
Formal Grammar

Languages & Grammars

Phrase-Structure Grammars

Types of Phrase-Structure Grammars

Derivation Trees

Backus-Naur Form

Intro to Languages

English grammar tells us if a given combination of words is a valid sentence.

The **syntax** of a sentence concerns its **form** while the **semantics** concerns its **meaning**.

e.g. the mouse wrote a poem

From a **syntax** point of view this is a valid sentence.

From a **semantics** point of view not so fast...perhaps in Disney land

Natural languages (English, French, Portuguese, etc) have very complex rules of syntax and not necessarily well-defined.

Formal Language

Formal language – is specified by well-defined set of rules of syntax

We describe the sentences of a formal language using a grammar.

Two key questions:

- 1 - Is a combination of words a valid sentence in a formal language?
- 2 – How can we generate the valid sentences of a formal language?

Formal languages provide models for both natural languages and programming languages.

Grammars

A formal *grammar* G is any compact, precise mathematical definition of a language L .

- As opposed to just a raw listing of all of the language's legal sentences, or just examples of them.

A grammar implies an algorithm that would generate all legal sentences of the language.

- Often, it takes the form of a set of recursive definitions.

A popular way to specify a grammar recursively is to specify it as a *phrase-structure grammar*.

Grammars (Semi-formal)

Example: A grammar that generates a **subset of the English language**

$$\langle sentence \rangle \rightarrow \langle noun_phrase \rangle \langle predicate \rangle$$
$$\langle noun_phrase \rangle \rightarrow \langle article \rangle \langle noun \rangle$$
$$\langle predicate \rangle \rightarrow \langle verb \rangle$$

$\langle \textit{article} \rangle \rightarrow a$

$\langle \textit{article} \rangle \rightarrow the$

$\langle \textit{noun} \rangle \rightarrow boy$

$\langle \textit{noun} \rangle \rightarrow dog$

$\langle \textit{verb} \rangle \rightarrow runs$

$\langle \textit{verb} \rangle \rightarrow sleeps$

A derivation of “the boy sleeps”:

$$\begin{aligned}\langle sentence \rangle &\Rightarrow \langle noun_phrase \rangle \langle predicate \rangle \\ &\Rightarrow \langle noun_phrase \rangle \langle verb \rangle \\ &\Rightarrow \langle article \rangle \langle noun \rangle \langle verb \rangle \\ &\Rightarrow the \langle noun \rangle \langle verb \rangle \\ &\Rightarrow the \ boy \langle verb \rangle \\ &\Rightarrow the \ boy \ sleeps\end{aligned}$$

A derivation of “a dog runs”:

$\langle sentence \rangle \Rightarrow \langle noun_phrase \rangle \langle predicate \rangle$
 $\Rightarrow \langle noun_phrase \rangle \langle verb \rangle$
 $\Rightarrow \langle article \rangle \langle noun \rangle \langle verb \rangle$
 $\Rightarrow a \langle noun \rangle \langle verb \rangle$
 $\Rightarrow a \text{ dog } \langle verb \rangle$
 $\Rightarrow a \text{ dog runs}$

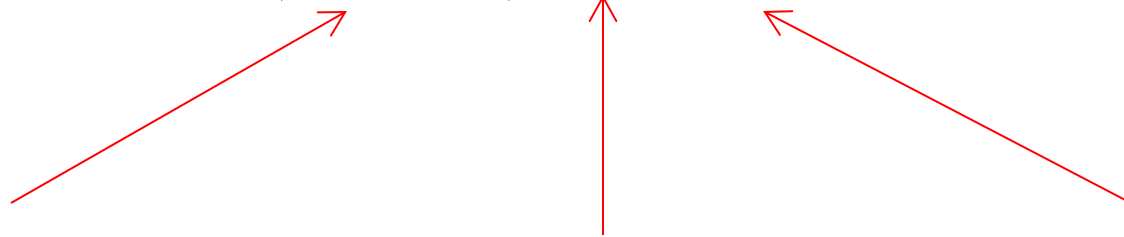
Language of the grammar:

$$L = \{ \text{“a boy runs”}, \\ \text{“a boy sleeps”}, \\ \text{“the boy runs”}, \\ \text{“the boy sleeps”}, \\ \text{“a dog runs”}, \\ \text{“a dog sleeps”}, \\ \text{“the dog runs”}, \\ \text{“the dog sleeps”} \}$$

Notation

$\langle noun \rangle \rightarrow boy$

$\langle noun \rangle \rightarrow dog$



Variable
or
Non-terminal

Production
rule

Terminal
Symbols of
the vocabulary

Symbols of
the vocabulary

Basic Terminology

- ▶ A ***vocabulary/alphabet***, V is a finite nonempty set of elements called symbols.
 - Example: $V = \{a, b, c, A, B, C, S\}$
- ▶ A ***word/sentence*** over V is a string of finite length of elements of V .
 - Example: Aba
- ▶ The ***empty/null string***, λ is the string with no symbols.
- ▶ V^* is the set of all words over V .
 - Example: $V^* = \{Aba, BBa, bAA, cab \dots\}$
- ▶ A ***language*** over V is a subset of V^* .
 - We can give some criteria for a word to be in a language.

Phrase-Structure Grammars

A *phrase-structure grammar* (abbr. PSG)

$G = (V, T, S, P)$ is a 4-tuple, in which:

- V is a vocabulary (set of symbols)
 - The “template vocabulary” of the language.
- $T \subseteq V$ is a set of symbols called *terminals*
 - Actual symbols of the language.
 - Also, $N \equiv V - T$ is a set of special “symbols” called *nonterminals*. (Representing concepts like “noun”)
- $S \in N$ is a special nonterminal, the *start symbol*.
 - in our example the start symbol was “sentence”.
- P is a set of *productions* (to be defined).
 - Rules for substituting one sentence fragment for another
 - Every production rule must contain at **least one nonterminal** on its left side.

Phrase-structure Grammar

► EXAMPLE:

- ❑ Let $G = (V, T, S, P)$,
- ❑ where $V = \{a, b, A, B, S\}$
- ❑ $T = \{a, b\}$,
- ❑ S is a start symbol
- ❑ $P = \{S \rightarrow ABa, A \rightarrow BB, B \rightarrow ab, A \rightarrow Bb\}$.

G is a **Phrase-Structure Grammar**.

What sentences can be generated
with this grammar?

Derivation

Definition

Let $G=(V,T,S,P)$ be a phrase-structure grammar.

Let $w_0=lz_0r$ (the concatenation of l , z_0 , and r) $w_1=lz_1r$ be strings over V .

If $z_0 \rightarrow z_1$ is a production of G we say that w_1 is **directly derivable** from w_0 and we write $w_0 \Rightarrow w_1$.

If w_0, w_1, \dots, w_n are strings over V such that $w_0 \Rightarrow w_1, w_1 \Rightarrow w_2, \dots, w_{n-1} \Rightarrow w_n$, then we say that w_n is derivable from w_0 , and write $w_0 \Rightarrow^* w_n$.

The sequence of steps used to obtain w_n from w_0 is called a **derivation**.

Language

Let $G(V, T, S, P)$ be a phrase-structure grammar. The language generated by G (or the language of G) denoted by $L(G)$, is the set of all strings of terminals that are derivable from the starting state S .

$$L(G) = \{w \in T^* \mid S \Rightarrow^* w\}$$

Language $L(G)$

► EXAMPLE:

Let $G = (V, T, S, P)$, where $V = \{a, b, A, S\}$, $T = \{a, b\}$, S is a start symbol and $P = \{S \rightarrow aA, S \rightarrow b, A \rightarrow aa\}$.

The language of this grammar is given by $L(G) = \{b, aaa\}$;

1. we can derive aA from using $S \rightarrow aA$, and then derive aaa using $A \rightarrow aa$.
2. We can also derive b using $S \rightarrow b$.

Another example

Grammar: $G=(V,T,S,P)$ $T=\{a,b\}$ $P =$

$$\begin{array}{l} S \rightarrow aSb \\ S \rightarrow \lambda \end{array}$$

$V=\{a,b,S\}$

Derivation of sentence $:ab$

$$S \Rightarrow aSb \Rightarrow ab$$

$S \rightarrow aSb$ $S \rightarrow \lambda$

$$S \rightarrow aSb$$

Grammar:

$$S \rightarrow \lambda$$

Derivation of sentence $aabb$

$$S \Rightarrow aSb \Rightarrow aaSbb \Rightarrow aabb$$



$$S \rightarrow aSb$$



$$S \rightarrow \lambda$$

Other derivations:

$$S \Rightarrow aSb \Rightarrow aaSbb \Rightarrow aaaSbbb \Rightarrow aaabbbb$$

$$S \Rightarrow aSb \Rightarrow aaSbb \Rightarrow aaaSbbb$$

$$\Rightarrow aaaaSbbbb \Rightarrow aaabbbbb$$

So, what's the language of the
grammar with the productions?

$$S \rightarrow aSb$$

$$S \rightarrow \lambda$$

Language of the grammar with the productions:

$$S \rightarrow aSb$$

$$S \rightarrow \lambda$$

$$L = \{a^n b^n : n \geq 0\}$$

PSG Example – English Fragment

We have $G = (V, T, S, P)$, where:

$V = \{(\text{sentence}), (\text{noun phrase}),$
 $(\text{verb phrase}), (\text{article}), (\text{adjective}),$
 $(\text{noun}), (\text{verb}), (\text{adverb}), a, the, large,$
 $hungry, rabbit, mathematician, eats, hops,$
 $quickly, wildly\}$

$T = \{a, the, large, hungry, rabbit, mathematician,$
 $eats, hops, quickly, wildly\}$

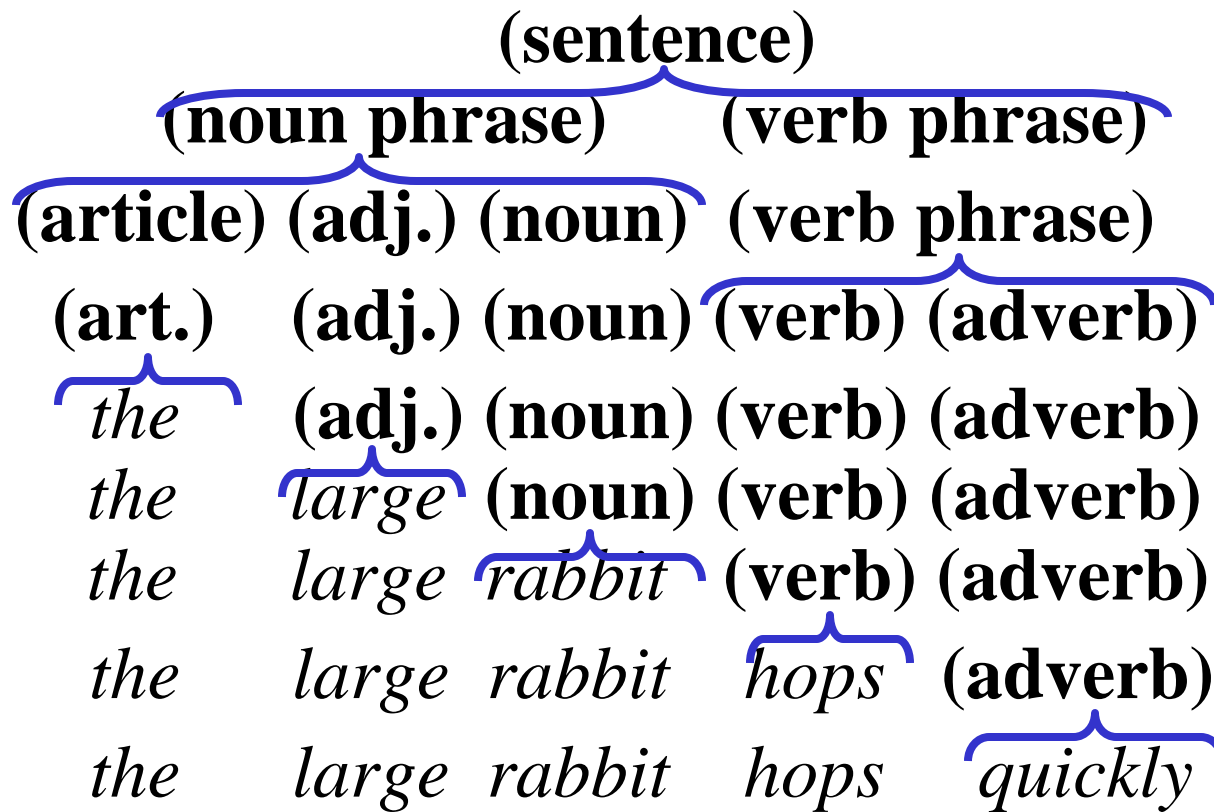
$S = (\text{sentence})$

$P = (\text{see next slide})$

Productions for our Language

$P = \{$ (sentence) \rightarrow (noun phrase) (verb phrase),
 (noun phrase) \rightarrow (article) (adjective) (noun),
 (noun phrase) \rightarrow (article) (noun),
 (verb phrase) \rightarrow (verb) (adverb),
 (verb phrase) \rightarrow (verb),
 (article) $\rightarrow a$, (article) $\rightarrow the$,
 (adjective) $\rightarrow large$, (adjective) $\rightarrow hungry$,
 (noun) $\rightarrow rabbit$, (noun) $\rightarrow mathematician$,
 (verb) $\rightarrow eats$, (verb) $\rightarrow hops$,
 (adverb) $\rightarrow quickly$, (adverb) $\rightarrow wildly \}$

A Sample Sentence Derivation



On each step, we apply a production to a fragment of the previous sentence template to get a new sentence template. Finally, we end up with a sequence of terminals (real words), **that is, a sentence of our language *L*.**

Another Example

Let $G = (\overbrace{\{a, b, A, B, S\}}^V, \overbrace{\{a, b\}}^T, S, \underbrace{P}_P)$.
 $\{S \rightarrow ABa, A \rightarrow BB, B \rightarrow ab, AB \rightarrow b\}$.

One possible derivation in this grammar is:

$$S \Rightarrow ABa \Rightarrow Aaba \Rightarrow BBaba \Rightarrow Bababa \\ \Rightarrow abababa.$$

Defining the PSG Types

Type 0: Phase-structure grammars – no restrictions on the production rules

Type 1: Context-Sensitive PSG:

- All after fragments are either longer than the corresponding before fragments, or empty:

$$\text{if } b \rightarrow a, \text{ then } |b| < |a| \quad \vee \quad a = \lambda .$$

Type 2: Context-Free PSG:

- All before fragments have length 1 and are nonterminals:

$$\text{if } b \rightarrow a, \text{ then } |b| = 1 \ (b \in N).$$

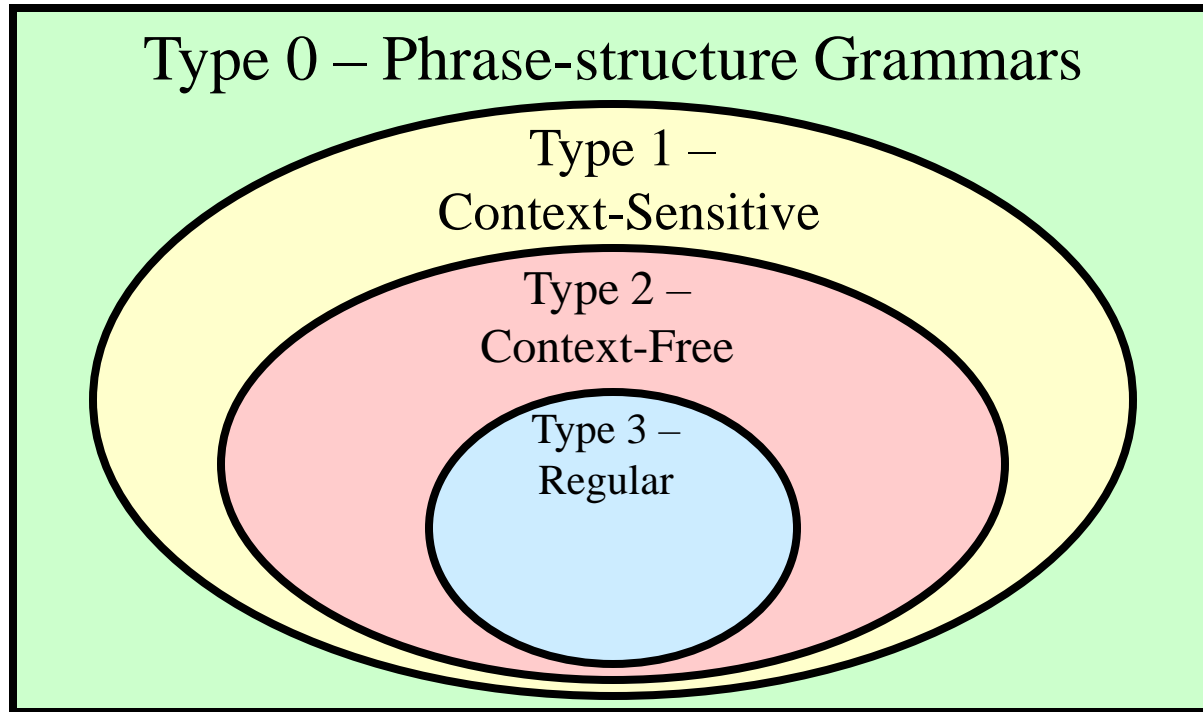
Type 3: Regular PSGs:

- All before fragments have length 1 and nonterminals
- All after fragments are either single terminals, or a pair of a terminal followed by a nonterminal.

$$\text{if } b \rightarrow a, \text{ then } a \in T \quad \vee \quad a \in TN.$$

Types of Grammars - Chomsky hierarchy of languages

Venn Diagram of Grammar Types:



Classifying grammars

Given a grammar, we need to be able to find the smallest class in which it belongs. This can be determined by answering three questions:

Are the left hand sides of all of the productions single non-terminals?

If yes, does each of the productions create at most one non-terminal and is it on the right?

Yes – regular

No – context-free

If not, can any of the rules reduce the length of a string of terminals and non-terminals?

Yes – unrestricted

No – context-sensitive

Definition: Context-Free Grammars

Grammar $G = (V, T, S, P)$

Vocabulary Terminal symbols Start variable

Productions of the form:

$A \rightarrow x$

Non-Terminal String of variables and terminals

Derivation Tree of A Context-free Grammar

-
- ▶ Represents the language using an ordered rooted tree.
 - ▶ **Root** represents the **starting symbol**.
 - ▶ **Internal vertices** represent the **nonterminal symbol** that arise in the production.
 - ▶ **Leaves** represent the **terminal symbols**.
 - ▶ If the production $A \rightarrow w$ arise in the derivation, where w is a word, the vertex that represents A has as children vertices that represent each symbol in w , in order from left to right.

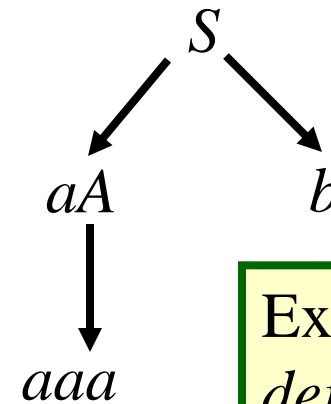
Language Generated by a Grammar

Example: Let $G = (\{S, A, a, b\}, \{a, b\}, S, \{S \rightarrow aA, S \rightarrow b, A \rightarrow aa\})$. What is $L(G)$?

Easy: We can just draw a tree of all possible derivations.

- We have: $S \Rightarrow aA \Rightarrow aaa$.
- and $S \Rightarrow b$.

Answer: $L = \{aaa, b\}$.



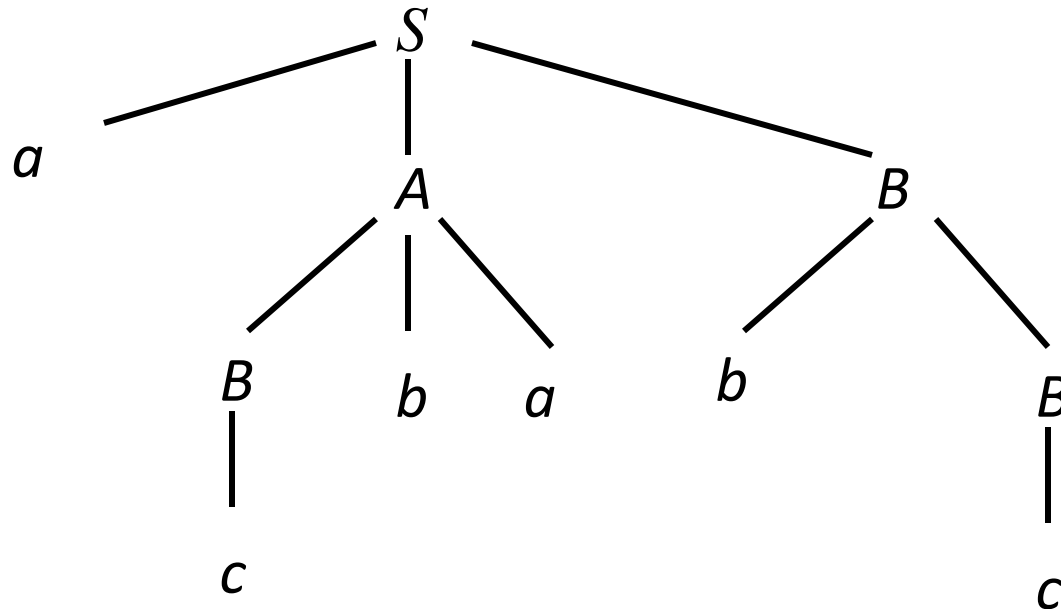
Example of a *derivation tree* or *parse tree* or *sentence diagram*.

Example: Derivation Tree

- Let G be a context-free grammar with the productions $P = \{S \rightarrow aAB, A \rightarrow Bba, B \rightarrow bB, B \rightarrow c\}$. The word $w = acbabc$ can be derived from S as follows:

$$S \Rightarrow aAB \rightarrow a(Bba)B \Rightarrow acbaB \Rightarrow acba(bB) \Rightarrow acbabc$$

Thus, the derivation tree is given as follows:



Backus-Naur Form

$\langle \text{sentence} \rangle ::= \langle \text{noun phrase} \rangle \langle \text{verb phrase} \rangle$

$\langle \text{noun phrase} \rangle ::= \langle \text{article} \rangle [\langle \text{adjective} \rangle] \langle \text{noun} \rangle$

$\langle \text{verb phrase} \rangle ::= \langle \text{verb} \rangle [\langle \text{adverb} \rangle]$

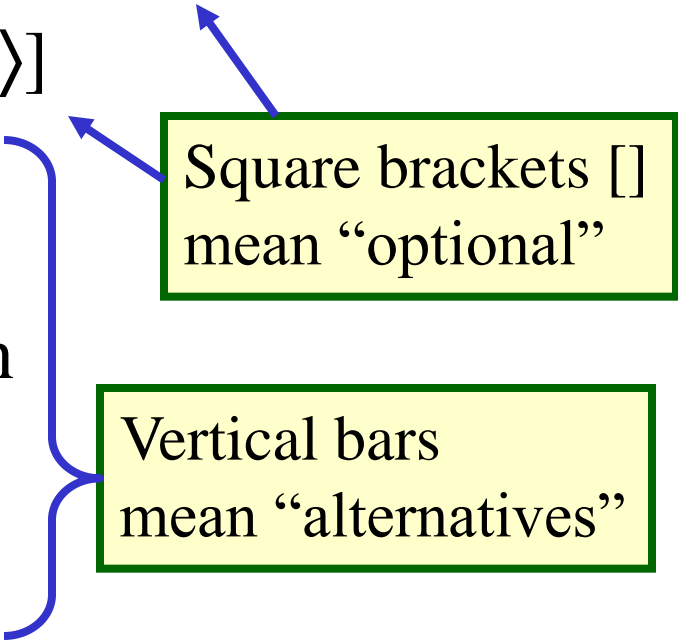
$\langle \text{article} \rangle ::= \text{a} / \text{the}$

$\langle \text{adjective} \rangle ::= \text{large} | \text{hungry}$

$\langle \text{noun} \rangle ::= \text{rabbit} | \text{mathematician}$

$\langle \text{verb} \rangle ::= \text{eats} / \text{hops}$

$\langle \text{adverb} \rangle ::= \text{quickly} / \text{wildly}$



Square brackets []
mean “optional”

Vertical bars
mean “alternatives”

Generating Infinite Languages

A simple PSG can easily generate an infinite language.

Example: $S \rightarrow 11S, S \rightarrow 0$ ($T = \{0,1\}$).

The derivations are:

- $S \Rightarrow 0$
- $S \Rightarrow 11S \Rightarrow 110$
- $S \Rightarrow 11S \Rightarrow 1111S \Rightarrow 11110$
- and so on...

$L = \{(11)^*0\}$ – the set of all strings consisting of some number of concatenations of 11 with itself, followed by 0.

Another example

Construct a PSG that generates the language $L = \{0^n 1^n \mid n \in \mathbf{N}\}$.

- 0 and 1 here represent symbols being concatenated n times, not integers being raised to the n th power.

Solution strategy: Each step of the derivation should preserve the invariant that the number of 0 's = the number of 1 's in the template so far, and all 0 's come before all 1 's.

Solution: $S \rightarrow 0S1, S \rightarrow \lambda$.

Context-Sensitive Languages

The language $\{ a^n b^n c^n \mid n \geq 1 \}$ is context-sensitive but not context free.

A grammar for this language is given by:

$$S \rightarrow aSBC \mid aBC$$

$$CB \rightarrow BC$$

$$aB \rightarrow ab$$

$$bB \rightarrow bb$$

$$bC \rightarrow bc$$

$$cC \rightarrow cc$$

Terminal
and
non-terminal



A derivation from this grammar is:-

$$\begin{aligned} S &\Rightarrow aSBC \\ &\Rightarrow aaBCBC && \text{(using } S \rightarrow aBC) \\ &\Rightarrow aabCBC && \text{(using } aB \rightarrow ab) \\ &\Rightarrow aabBCC && \text{(using } CB \rightarrow BC) \\ &\Rightarrow aabbCC && \text{(using } bB \rightarrow bb) \\ &\Rightarrow aabbcC && \text{(using } bC \rightarrow bc) \\ &\Rightarrow aabbcc && \text{(using } cC \rightarrow cc) \end{aligned}$$

which derives $a^2b^2c^2$.