

Lecture 11

Moving From Strings to Trees

We have been able to show that English is not a regular language due to the structural complexity of center embedding. This leaves us with two choices: we may either stick with finite-state models and claim that center embedding does not pose a challenge in practice, or we can once again move to a more expressive model. The first choice is a valid option in a variety of applications but requires some major sacrifices. We will see that these sacrifices cannot be reconciled with our project of exploring the properties of language from a computational perspective, and thus we will take a hint from linguists and expand our model by moving from strings to trees.

1 Finite-State Methods and the Role of Unboundedness

As mentioned last time, the reason that center-embedding patterns — abstractly represented as $a^n b^n$ — are beyond the capabilities of finite-state methods is that there is no upper bound on the depth of embedding. An FSA has to memorize the exact number of a s in order to ensure that the same number of b s follow, but since an FSA has a fixed number k of states, it can only partition strings into $k + 1$ distinct equivalence classes (one equivalence class for each state, plus one for strings for which there is no defined transition). Hence some strings with distinct numbers of a s must wind up in the same equivalence class and thus can be followed by the same number of b s according to the automaton, which is clearly not the case for the language $a^n b^n$.

This argument is mathematically flawless, but it faces a major empirical challenge: center-embedding is **not** unbounded in the data we have access to. No speaker spontaneously produces a sentence with ten levels of center-embedding, and even if we found a speaker with a propensity for convoluted center embedding constructions and convinced that person to spend the rest of their life uttering a single sentence full of center embeddings, that person could only produce a fixed number of embeddings before they die of old age. Even if we somehow could found a tradition of center-embedding performance art, where one speaker starts a sentence with center embeddings that are then continued by another speaker, and then another, on and on until the end of time, we could only reach a final level of embeddings before the universe collapses in on itself. We live in a finitistic universe, wherefore unboundedness is not an empirically verifiable property.

A slightly different, more application-oriented argument posits that even if humans might be theoretically capable of unbounded center-embedding, there is little reason to incorporate this property into our model if they never make use of it. Why make your

The ideal solution would be an elegant tool that can be automatically translated into an efficient one if necessary, but that is not always feasible.

model more powerful if that power is never needed? This is indeed a valid concern for industrial-grade applications, where time equals money and programs should run as quickly as possible. But even there this issue is not cut and dry. After all, “time equals money” also means that programs should be easy to extend, modify and maintain, since the manhours spent on these tasks do not come for free. When given a choice between an efficient but complicated tool on the one hand or a slower yet elegant tool on the other, the latter might be the better solution. In addition, the elegant tool might actually be the faster one in practice — just because it can theoretically perform much worse does not entail that it does so for the specific task at hand. As you can see, the unboundedness question does not have a clear cut answer, it all depends on what your goals are.

Our inquiry is driven by scientific curiosity, so questions of efficiency affect us only to the extent that the resource usage of our model has to be reconcilable with what we know about human cognition. And even this restriction does not outweigh our desire to state insightful generalizations. That is why we put such a high premium on abstraction: by deliberately excluding certain factors we can home in on broad, appealing generalizations. These generalizations still hold in a more detailed model that is closer to the wetware, but they are much harder to discern due to all the complicating factors that come with a more faithful model.

For our purposes, unboundedness is an essential assumption because boundedness acts as a great equalizer that pushes everything within the bounds of finite-state machinery. Recall that we assumed in our discussion of phonology that some long-distance processes are unbounded even though for all practical purposes the length of words is bounded in the same way as the number of center embeddings. We did this because it allowed us to formalize important differences between local and non-local processes without losing track of their commonalities. Similarly, assuming that center embedding is unbounded brings out an important difference between this non-regular kind of embedding and right embedding, which is still regular. This contrast is even reflected in human processing, where right embeddings are much easier to parse than center embeddings. Distinguishing bounded center embedding from bounded right embedding would be more involved.

A brief glance at an FSA for bounded center embedding also reveals that redundancy of this account. Each level of embedding corresponds to a specific subgraph of the automaton, but all these subgraphs look exactly the same. So we have a big number of states that all do exactly the same work, the only difference is that one is used in embedding level 2 and another one in level 5. This also raises the question why natural language grammars treat all those levels exactly the same — if the automaton distinguishes level 2 from level 3, why can't 3 use the opposite word order of 2? This level-sensitive automaton would have exactly the same number of states as the one that treats all levels of embedding the same. If we assume that center embedding is unbounded, though, then it might be possible to show that level-sensitive formalisms are more complicated than those that treat all levels the same.

To sum up, unboundedness is not an undisputable fact, first and foremost it is yet another abstraction in the service of exploring specific questions. However, what we find may serve as indirect evidence for unboundedness when embedded in a system of ancillary assumptions. For example, the enormous size of FSAs with bounded center embedding and the lack of level-sensitive embeddings in natural language suggest that unboundedness is cognitively real, at least if one assumes that small, succinct grammars are preferred for some reason. With other constructions, there may be less

of a reason to assume that they are unbounded (for instance multiple wh-movement), so we will always have to weigh carefully what the benefits of unboundedness may be on a case-by-case basis.

2 Tree Languages

2.1 Context-Free Grammars

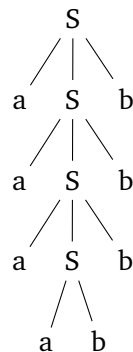
While unbounded center embedding exceeds the limits of finite-state methods, it is easily accounted for with phrase structure rules. The language $a^n b^n$, $n \geq 1$, is generated by two rules, which can even be conflated into a single one with some syntactic sugar:

$$S \rightarrow ab \mid aSb$$

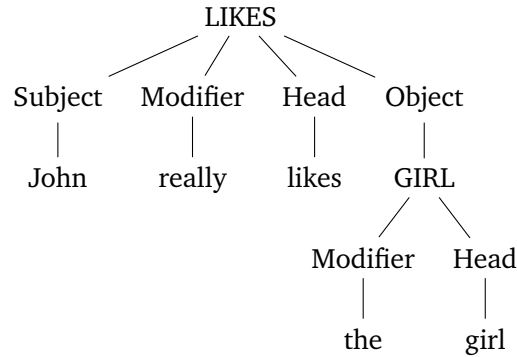
We can interpret this rule as a mechanism for rewriting strings, where we start with S and then apply rules until no more rewriting steps are possible.

Rule	String
start	S
$S \rightarrow aSb$	aSb
$S \rightarrow aSb$	$aaSbb$
$S \rightarrow aSb$	$aaaSbbb$
$S \rightarrow ab$	$aaaabbbb$

Linguists are more familiar with the tree-based representation of rule application.



Phrase structure grammars are also known as *context-free grammars* (CFGs). The latter term focuses on the rule format, which must be of the form $A \rightarrow \alpha$, where α is a string of symbols. These rules lack the context specification used by SPE, among others, so they are indeed context-free. The term phrase structure grammar instead focuses on what the grammar is meant to describe, namely the phrase structure of sentences. Obviously CFGs can be used to describe other kinds of structures. For instance, the sentence *John really likes the girl* could be assigned the tree below, which represents the functional relations between words (this tree is inspired by Dependency Grammar, which we will discuss at a later point).



As you can see, the term context-free grammars is slightly more general even though it refers to exactly the same kind of mathematical object. Fans of generality that we are, we will henceforth speak of CFGs rather than phrase structure grammars.

Definition 11.1 (CFG). A context-free grammar is a triple $G := \langle \Sigma, S, R \rangle$, where

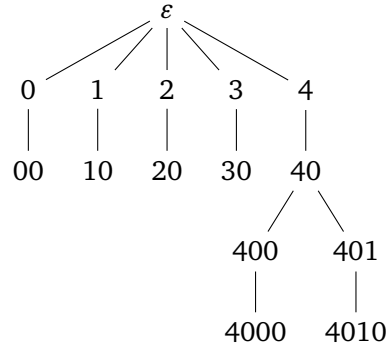
- Σ is an alphabet,
- $S \in \Sigma$ is the *start symbol*,
- R is finite set of rules $A \rightarrow \alpha$ such that $A \in \Sigma$ and $\alpha \in \Sigma^*$.

A symbol $a \in \Sigma$ is *terminal* iff R contains no rule of the form $a \rightarrow \alpha$. Otherwise a is *non-terminal*. The corresponding subsets of Σ are denoted T_Σ and N_Σ . We always require $S \in N_\Sigma$. The language $L(G)$ generated by G is the smallest set containing all strings that 1) can be obtained from S by finitely many applications of rules in R , and 2) contain no non-terminal.

Since CFGs can handle center embedding and are equivalent to the familiar linguistic formalism of phrase structure rules, they are a promising starting point for a formal model of syntax. We will soon see that they are indeed very closely related to a model that we have explored in great detail.

2.2 Trees and Tree Languages

Before we move on, it will be useful to formalize trees. This can be done in a plethora of ways, but we will opt for the definition in terms of *Gorn domains* (Gorn 1967) because it couples each node with an address that directly represents its location in the tree. Intuitively, a Gorn domain is a set of addresses of the form $m \cdot l$, where m is the address of the node's mother and l the number of left siblings it has. So the tree for *John really likes the girl* that was given in the previous section would be associated with node addresses as shown below:



The address for the root is ε since it has neither a mother nor a left sibling. For all other nodes, the formula $m \cdot l$ applies as described.

Notice that an entry like 40 is read “four-zero” since it is the leftmost daughter of the fifth daughter of the root, whereas 40 “forty” would refer to the 41st daughter of the node. This ambiguity is due to the decimal system being incapable of representing the difference between $4 \cdot 0$ and 40. Strictly speaking, an address like 401 should actually be written as $4 - 0 - 1$ to distinguish it from $40 - 1$ and 401, but this creates clutter that is best avoided when possible.

Definition 11.2 (Gorn domain). A Gorn domain D is a subset of \mathbb{N}^* such that

dominance closure $ui \in D$ implies $u \in D$,

left sibling closure $ui \in D$ implies $uj \in D$ for all $0 \leq j < i$.

We call $u \in D$ a *leaf* iff there is no $i \in \mathbb{N}$ such that $ui \in D$. Given some subset S of D , ui is a *root* of S iff $u \notin S$.

A tree is a Gorn domain that maps each node address to a node label.

Definition 11.3 (Tree). A (finite) Σ -tree is a pair $t := \langle D, \ell \rangle$, where

- D is a (finite) Gorn Domain,
- $\ell : D \rightarrow \Sigma$ is a total function.

The *string yield* $yd(t)$ of t is the longest string $s_1 \cdots s_n$ such that

- s_i is a leaf of D ($0 \leq i \leq n$),
- for $s_i := mu$ and $s_j := nv$ ($m, n \in \mathbb{N}$, $u, v \in \mathbb{N}^*$), $i < j$ iff $m < n$.

The *depth* of subtree t with root u is the length of the longest i such that $ui \in D$.

Unless indicated otherwise, all trees are assumed to be finite.

Since we now have both strings and trees as mathematical objects, it makes sense to distinguish between *string languages* and *tree languages*. The former are sets of strings, the latter sets of trees. All the languages we have seen so far were string languages.

Definition 11.4 (Tree Language). A tree language L over Σ is a set of Σ -trees. Its string yield is the string language $yd(L) := \{yd(t) \mid t \in L\}$.

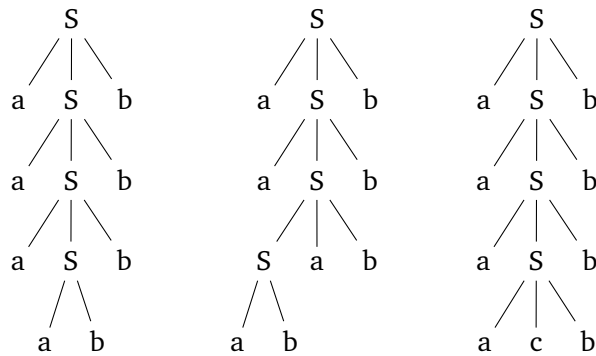
2.3 Tree Languages of Context-Free Grammars

Our definition of CFGs defines the string language generated by the grammar, but not the tree language even though we saw that sequences of rewriting steps can be represented as trees. The intuition for going from rewriting rules to trees is simple enough:

1. Draw a node with label S .
2. If A is rewritten as $\alpha := a_1 \cdots a_n$, add a_1, \dots, a_n as daughters of the corresponding node (in the same left-to-right order).

This way a CFG can be treated as a mechanism for building trees rather than just strings. But what exactly is the tree language generated by some random CFG?

Let us first think about whether certain trees are generated by a specific grammar, and why this is the case. Below you have several trees. The left one is generated by the CFG $\langle \{S, a, b\}, \{S \rightarrow aSb, S \rightarrow ab\} \rangle$, the others are not.



The left one satisfies all the conditions the rewrite rules establish with respect to the mother-of and the left-sibling relation. The tree in the middle contains a subtree where S is a left sibling of a , which can never happen with the specified rewrite rules. The third tree contains c as a daughter of S , which is also impossible. Notice that all these violations can be verified in a local fashion: we only have to look at a node and its daughters to find ill-formed subtrees. So we can postulate the following:

Definition 11.5 (CFG Tree Language). A CFG G generates a tree t iff

- the root of t is the start symbol, and
- all leaves of t are terminal, and
- for every subtree s of t such that s is of the form $[_A A_1 \cdots A_n]$, G contains a rule $A \rightarrow A_1 \cdots A_n$.

The tree language of G is the set of all trees that are generated by G .

This is a definition, not a theorem. We are simply describing what kind of trees are built via the translation procedure used by linguists. But we can show that this definition is useful in the sense that the tree language yields the same string language as the grammar.

Theorem 11.6. Let G be a CFG that generates string language L and tree language T . Then $L = \text{yd}(T)$. \lrcorner

Proof. To see that $L \subseteq \text{yd}(T)$, take any string $w \in L$. By definition, w was obtained from S via a finite number of applications of rewrite rules of G . The standard translation from such rule applications to trees yields a tree that is generated by G and has w as its string yield.

In the other direction $\text{yd}(T) \subseteq L$ follows from the fact that if t is generated by G , then each subtree $[_A A_1 \cdots A_n]$ is matched by a rewrite rule of G . Hence there is a sequences of rewrite rules that produces the string yield of t . Since t 's string yield consists only of terminal symbols, and t 's root is S , $\text{yd}(t)$ is generated by G . \square

2.4 Strictly Local Tree Languages

The condition that every subtree of depth 1 must be matched by a rewrite rule could be simplified by converting each rewrite rule into a tree of depth 1. Instead of a finite set of rewrite rules, one would then have a finite set of subtrees of depth 1, and a tree is generated by the grammar iff all its subtrees of depth 1 are included in this set. This idea should sound awfully familiar to you: it is the tree analogue of strictly local grammars.

Definition 11.7 (k -trees). A k -tree over alphabet Σ is a tree of depth $k - 1$. A k -augment is a unary branching tree of depth $k - 1$ where every node is labeled \times . Given a Σ -tree t , its k -augmented counterpart \hat{t}_k is the result of adding a k -augment above the root and below each leaf. The set of k -trees of a tree t is given by $2\text{-trees}(t) := \{s \mid s \text{ is a subtree of } t \text{ with depth } k - 1\}$.

Definition 11.8 (Strictly Local Tree Language). A finite set of k -trees is called a *strictly k -local tree grammar*. A positive strictly k -local tree grammar G generates the tree language $L(G) := \{t \mid 2\text{-trees}(t) \subseteq G\}$. A negative strictly k -local tree grammar G generates the tree language $L(G) := \{t \mid 2\text{-trees}(t) \cap G = \emptyset\}$. A tree language L is strictly k -local iff it is generated by some strictly k -local tree grammar. The class SL^T of *strictly local tree languages* is given by $\bigcup_{k \geq 1} \{L \mid L \text{ is strictly } k\text{-local}\}$.

Lemma 11.9. For every positive strictly k -local grammar, there is a negative strictly k -local grammar that generates the same tree language, and *vice versa*. \lrcorner

Proof. A simple extension of the proof for strictly k -local string grammars. \square

Theorem 11.10. The class of tree languages generated by CFGs is properly included in the class of strictly 2-local tree languages. \lrcorner

Proof. For every CFG G one can construct an equivalent strictly 2-local tree grammar G_2 . For every rewrite rule $A \rightarrow A_1 \cdots A_n$ of G , G_2 contains the 2-tree $[_A A_1 \cdots A_n]$. For every terminal symbol a , we add $[_a \times]$. Finally, G also contains $[_\times S]$. It is easy to see that G_2 generates exactly the same tree language.

Inclusion is proper because a symbol can be both terminal and non-terminal in a strictly 2-local tree language. Consider for instance the grammar $\{[_\times S], [_S \times]\}$, which generates only the tree S . \square

Theorem 11.11. For all $k \geq 1$, $SL_{k-1}^T \subsetneq SL_k^T$. ┘

Theorem 11.12. For all $k \geq 1$, $yd(SL_{k-1}^T) = yd(SL_k^T)$. ┘