

Syntax: Introduction

Till now, we have discussed:

Language modeling — understanding ordering of words

POS tagging — understanding roles of words

Now we will discuss more complex notions, such as grammatical relations among words

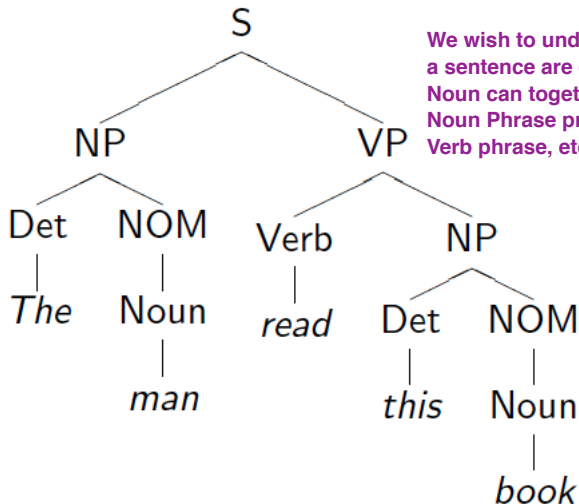
What is Syntax?

- Refers to the way words are arranged together, and the relationship between them.

Syntax Tree: Example

Sentence: The man read this book

We wish to understand how different words in a sentence are grouped, e.g., a Det and a Noun can together give a Noun Phrase, a Noun Phrase preceded by a verb can give a Verb phrase, etc.



Our final goal:
How to construct a syntax
tree for a given sentence?

Defining the notions: Constituency

Constituent

A group of words acts as a single unit - phrases, clauses etc.

Defining the notions: Constituency

Constituent

A group of words acts as a single unit - phrases, clauses etc.

Part of Speech - "Substitution Test"

The {sad, intelligent, green, fat, ...} one is in the corner.

Different words having the same POS (here, adjective) can be used as a substitute for each other.

Defining the notions: Constituency

Constituent

A group of words acts as a single unit - phrases, clauses etc.

Part of Speech - "Substitution Test"

The {sad, intelligent, green, fat, ...} one is in the corner.

Constituency: Noun Phrase

- *Kermit the frog*
- *they*
- *December twenty-sixth*
- *the reason he is running for president*

Similarly, different word-sequences (phrases) can be used as a substitute for each other.
E.g., noun phrases

Constituent Phrases

How are constituent phrases usually named?

Usually named based on the word that heads the constituent:

<i>the man from Amherst</i>	is a Noun Phrase (NP) because the head <i>man</i> is a noun
<i>extremely clever</i>	is an Adjective Phrase (AP) because the head <i>clever</i> is an adjective
<i>down the river</i>	is a Prepositional Phrase (PP) because the head <i>down</i> is a preposition
<i>killed the rabbit</i>	is a Verb Phrase (VP) because the head <i>killed</i> is a verb

Constituent Phrases

Usually named based on the word that heads the constituent:

<i>the man from Amherst</i>	is a Noun Phrase (NP) because the head <i>man</i> is a noun
<i>extremely clever</i>	is an Adjective Phrase (AP) because the head <i>clever</i> is an adjective
<i>down the river</i>	is a Prepositional Phrase (PP) because the head <i>down</i> is a preposition
<i>killed the rabbit</i>	is a Verb Phrase (VP) because the head <i>killed</i> is a verb

Words can also act as phrases

Joe grew potatoes

Joe and *potatoes* are both nouns and noun phrases

Constituent Phrases

Usually named based on the word that heads the constituent:

<i>the man from Amherst</i>	is a Noun Phrase (NP) because the head <i>man</i> is a noun
<i>extremely clever</i>	is an Adjective Phrase (AP) because the head <i>clever</i> is an adjective
<i>down the river</i>	is a Prepositional Phrase (PP) because the head <i>down</i> is a preposition
<i>killed the rabbit</i>	is a Verb Phrase (VP) because the head <i>killed</i> is a verb

Words can also act as phrases

Joe grew potatoes

Joe and *potatoes* are both nouns and noun phrases

Compare with: *The man from Amherst grew beautiful russet potatoes.*

Joe appears in a place that a larger noun phrase could have been.

Evidence that constituency exists

They appear in similar environments

Evidence that constituency exists

They appear in similar environments

Kermit the frog comes on stage

They come to Massachusetts every summer

December twenty-sixth comes after Christmas

The reason he is running for president comes out only now.

But not each individual word in the constituent

*The comes out... *is comes out... *for comes out...

Evidence that constituency exists

They appear in similar environments

Kermit the frog comes on stage

They come to Massachusetts every summer

December twenty-sixth comes after Christmas

The reason he is running for president comes out only now.

But not each individual word in the constituent

*The comes out... *is comes out... *for comes out...

Can be placed in a number of different locations

Evidence that constituency exists

They appear in similar environments

Kermit the frog comes on stage

They come to Massachusetts every summer

December twenty-sixth comes after Christmas

The reason he is running for president comes out only now.

But not each individual word in the constituent

*The comes out... *is comes out... *for comes out...

Can be placed in a number of different locations

Constituent = Prepositional phrase: On December twenty-sixth

On December twenty-sixth I'd like to fly to Florida.

I'd like to fly on December twenty-sixth to Florida.

I'd like to fly to Florida on December twenty-sixth.

But not split apart

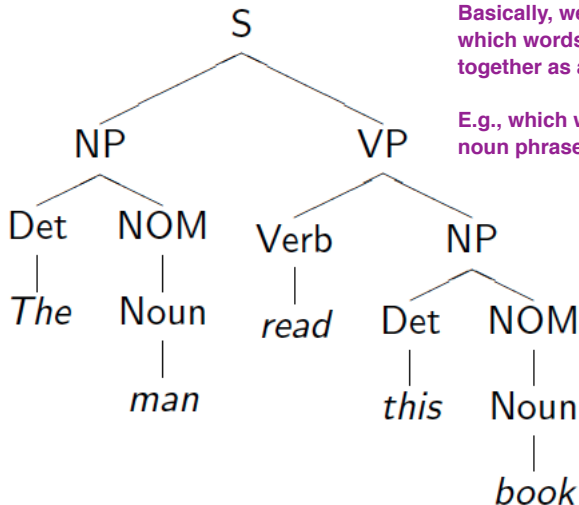
*On December I'd like to fly twenty-sixth to Florida.

*On I'd like to fly December twenty-sixth to Florida.

Modeling Constituency: what tool do we need?

Basically, we want tools to understand which words (in a given text) come together as a unit.

E.g., which words together make up a noun phrase?



Modeling Constituency

Context-free grammar

The most common way of modeling constituency

Modeling Constituency

Context-free grammar

The most common way of modeling constituency

Consists of production Rules

These rules express the ways in which the symbols of the language can be grouped and ordered together

Modeling Constituency

Context-free grammar

The most common way of modeling constituency

Consists of production Rules

These rules express the ways in which the symbols of the language can be grouped and ordered together

Example

Noun phrase can be composed of either a ProperNoun or a determiner (Det) followed by a Nominal; a Nominal can be more than one nouns

Modeling Constituency

Context-free grammar

The most common way of modeling constituency

Consists of production Rules

These rules express the ways in which the symbols of the language can be grouped and ordered together

Example

Noun phrase can be composed of either a ProperNoun or a determiner (Det) followed by a Nominal; a Nominal can be more than one nouns

$NP \rightarrow \text{Det Nominal}$

$NP \rightarrow \text{ProperNoun}$

$\text{Nominal} \rightarrow \text{Noun} \mid \text{Noun Nominal}$

CFG: $G = (T, N, S, R)$

- T : set of terminals
- N : set of non-terminals
 - ▶ For NLP, we distinguish out a set $P \subset N$ of pre-terminals, which always rewrite as terminals
- S : start symbol
- R : Rules/productions of the form $X \rightarrow \gamma$, $X \in N$ and $\gamma \in (T \cup N)^*$

CFG: $G = (T, N, S, R)$

- T : set of terminals
- N : set of non-terminals
 - ▶ For NLP, we distinguish out a set $P \subset N$ of pre-terminals, which always rewrite as terminals
- S : start symbol
- R : Rules/productions of the form $X \rightarrow \gamma$, $X \in N$ and $\gamma \in (T \cup N)^*$

Terminals and pre-terminals

Terminals mainly correspond to words in the language while pre-terminals mainly correspond to POS categories

Example

$NP \rightarrow \text{Det Nominal}$

$NP \rightarrow \text{ProperNoun}$

$\text{Nominal} \rightarrow \text{Noun} \mid \text{Noun Nominal}$

“Noun”, “ProperNoun”, “Det” are pre-terminals.

But we cannot generate any string with these rules, since there is no terminal.

Example

$NP \rightarrow \text{Det Nominal}$

$NP \rightarrow \text{ProperNoun}$

$\text{Nominal} \rightarrow \text{Noun} \mid \text{Noun Nominal}$

Now, these can be combined with other rules, that express facts about a lexicon.

Example

$NP \rightarrow \text{Det Nominal}$

$NP \rightarrow \text{ProperNoun}$

$\text{Nominal} \rightarrow \text{Noun} \mid \text{Noun Nominal}$

Now, these can be combined with other rules, that express facts about a lexicon.

$\text{Det} \rightarrow \text{a}$

$\text{Det} \rightarrow \text{the}$

$\text{Noun} \rightarrow \text{flight}$

Terminals: “a”, “the”, “flight”

Pre-terminals: “Det”, “Noun”, “ProperNoun”

Non-terminals: “NP”, “Nominal”

CFG as a generator

NP → Det Nominal

NP → ProperNoun

Nominal → Noun | Noun Nominal

Det → a

Det → the

Noun → flight

CFG as a generator

$NP \rightarrow \text{Det Nominal}$

$NP \rightarrow \text{ProperNoun}$

$\text{Nominal} \rightarrow \text{Noun} \mid \text{Noun Nominal}$

$\text{Det} \rightarrow \text{a}$

$\text{Det} \rightarrow \text{the}$

$\text{Noun} \rightarrow \text{flight}$

Generating 'a flight':

CFG as a generator

NP → Det Nominal

NP → ProperNoun

Nominal → Noun | Noun Nominal

Det → a

Det → the

Noun → flight

Generating 'a flight':

NP → Det Nominal

CFG as a generator

$NP \rightarrow \text{Det Nominal}$

$NP \rightarrow \text{ProperNoun}$

$\text{Nominal} \rightarrow \text{Noun} \mid \text{Noun Nominal}$

$\text{Det} \rightarrow \text{a}$

$\text{Det} \rightarrow \text{the}$

$\text{Noun} \rightarrow \text{flight}$

Generating 'a flight':

$NP \rightarrow \text{Det Nominal}$

$\rightarrow \text{Det Noun}$

CFG as a generator

$NP \rightarrow \text{Det Nominal}$

$NP \rightarrow \text{ProperNoun}$

$\text{Nominal} \rightarrow \text{Noun} \mid \text{Noun Nominal}$

$\text{Det} \rightarrow \text{a}$

$\text{Det} \rightarrow \text{the}$

$\text{Noun} \rightarrow \text{flight}$

Generating 'a flight':

$NP \rightarrow \text{Det Nominal}$

$\rightarrow \text{Det Noun} \rightarrow \text{a Noun}$

CFG as a generator

$NP \rightarrow \text{Det Nominal}$

$NP \rightarrow \text{ProperNoun}$

$\text{Nominal} \rightarrow \text{Noun} \mid \text{Noun Nominal}$

$\text{Det} \rightarrow \text{a}$

$\text{Det} \rightarrow \text{the}$

$\text{Noun} \rightarrow \text{flight}$

Generating 'a flight':

$NP \rightarrow \text{Det Nominal}$

$\rightarrow \text{Det Noun} \rightarrow \text{a Noun} \rightarrow \text{a flight}$

CFG as a generator

$NP \rightarrow \text{Det Nominal}$

$NP \rightarrow \text{ProperNoun}$

$\text{Nominal} \rightarrow \text{Noun} \mid \text{Noun Nominal}$

$\text{Det} \rightarrow \text{a}$

$\text{Det} \rightarrow \text{the}$

$\text{Noun} \rightarrow \text{flight}$

Generating 'a flight':

$NP \rightarrow \text{Det Nominal}$

$\rightarrow \text{Det Noun} \rightarrow \text{a Noun} \rightarrow \text{a flight}$

- Thus a CFG can be used to randomly generate a series of strings

CFGs and Recursion

A CFG can contain recursive rules.

E.g., the rule for a Prepositional Phrase (PP)

Recursive Definition

- $PP \rightarrow \text{Prep } NP$
- $NP \rightarrow \text{Noun } PP$

Recursive Definition

- $PP \rightarrow \text{Prep NP}$
- $NP \rightarrow \text{Noun PP}$

Example Sentence

[_SThe mailman ate his [_{NP} lunch [_{PP} with his friend [_{PP} from the cleaning staff [_{PP} of the building [_{PP} at the intersection [_{PP} on the north end [_{PP} of town]]]]]]].

start of a NP

start of a NP

start of a NP

start of a NP

start of a NP

A CFG defines a formal language = set of all sentences (string of words) that can be derived by the grammar

- Sentences in this set are said to be **grammatical**
- Sentences outside this set are said to be **ungrammatical**