

Natural Language Processing (NLP)

The Research field - essentially.
Gold Standard Algorithms - > English.
Natural language processing algorithms are particularly tested in English.

English syntax is very simple -> English maintains that a simple syntax indicates an evolved language. French is an example of an unevolved language. 72% -> indicates a very good result (success). We are carrying a syntactic analysis to a human subject. Morphologically, syntactically processing.

NLP refers to a section of

Trying to express natural language structures formally. We capture the difference between linguistics and language processing. This is NOT linguistics. We use language from the point of view of data scientists. Software has to not only understand language structure but also as a reverse process to generate natural language.

Natural Language Generation

list of topics for the exam.
let - applied system.

NLP - creating a technology that deals with natural language (entirely different from formal languages techniques). Computational linguistics tries to find models for natural languages. NLP takes place at the following levels: (Chapter 15 completely different).

- NLP Processing Levels:
- morphological
 - syntactical
 - semantic
 - pragmatic level.

Morphological → Input phrase, but not with the sense of a set of sentences! A phrase is any group of words!

I want to determine the part of speech of each word in a phrase.

Ultimate goal is for the software to understand a phrase. Natural language is characterized by ambiguity

(for ex: bank, banking ops.)

name verb

Ambiguity is very strong in English. Semantic ambiguity is particularly high in English.

As a noun → 7 meanings. 12 meanings.

As a verb → 5 meanings

At the lowest level, we perform morphological analysis → part of speech is v. important, it helps solve ambiguity.

[Tag each word in the input phrase with a part of speech. → Software is termed POS-tagger → family of software, Part of Speech taggers, there are several available, the

Heavens being the Stanford Group one.

(Stanford part of speech feger, famous and very successful)

At the PC, we will work with NLTRE, Pithon, which has the same PC-feger. It is not clear which one is the best. We don't have errors comparison yet.

The programmers always use the Stanford POS-feger.

Stanford POS-feger had 98% success which is extremely high. This is however a field where too much success cannot be obtained. At that point, it was considered that morphological feger was solved and was considered a closed field of NLP.

Syntactic level → means the structure of the

phrase (there are completely different syntactic levels)

Phrase structure

Syntactic level is very much established. Discussion is focused on the type of syntactic analysis.

- Semantic level → is very much alive, even for English, nothing is established.

Semantics → means determining meaning, sense.

At semantic level, I am trying to determine the meaning of an input phrase, context free,

regardless of the context it occurs in.

Various types of ambiguities that can occur.

Morphological ambiguity → Ree flies like sand
verb? noun?

Establishing meaning context free.

— At the pragmatic level, we do take

Context and context into account.

This is an underdeveloped level → many PhD theses in this field.

We will not get into the pragmatic level at all.

We will focus on the semantic level.

There are 2 completely different theories concerning syntax.

— Syntactic level focuses on the structure of a phrase. Determining the structure of a phrase → the process is called parsing.

① Syntactic analysis algorithms are called parsing algorithms.

Parsing is very well established for English.

Terms for syntactic analysis:

— syntactic phrase (group syntactic)

↳ a group of words which are organized around a word, which is considered the group head. Those words are in cohesion around that group head.

This selection is insured by grammar rules.

Various types of syntactic groups:

- noun phrase (NP) - the head is a noun.
- verb phrase (VP) - the head is a verb.

- adjective phrase (AdjP)

- adverbial phrase (AdvP)

- prepositional phrase (PP) - the head is a preposition.

Determining the structure of a phrase (syntactical analysis) can mean obtaining the concatenated phrases that represent the rest of the phrase.

It of view the structure of a phrase as a concatenation of well use in the background a generative grammar, which gives rules.

Rules: accordance $S \rightarrow NPVP$

$VP \rightarrow V PP$

VP phase \rightarrow verb
PP prep. phrase.

$PP \rightarrow PP \cdot NP$

These are recursive structure rules, applied for generative grammars (Chomsky).

If we change the type of grammar then the syntactic analysis changes \rightarrow the output changes.

Bibliografie → pentru analiză sintactică.

1. F. Hristea - Introducere în procesarea limbajului natural, Editura Universităţii din Bucureşti, 2010.
↳ partea de gramatică generativă este identică în cele 2 cărţi recomandate. (capitolul 3, intitulat "Gramatică").

2. Aspecte ale reprezentării cunoştinţelor, cu aplicaţii în procesarea limbajului natural.

Ed. Univ. Buc., 2011, capitolul 3 - "Gramatică generativă pt. procesarea limbajului natural".

Astea cărţi sunt pt. partea introductivă.

- Analiză sintactică bazată pe constituenţi: Constituantul este un grup sintactic.
- Analiză sintactică de dependenţă este complet diferită.

Tipurile de analiză sintactică → depind de tipul de aplicaţie.

Ultimate goal - make it easier for the software to understand natural language.

Interdisciplinarity.

This type of syntactic analysis uses generative grammars, which are tools of knowledge representation. A grammar

Generative grammars → set of phrase-structure rules, which are recursive rules based on concatenation (coined by Chomsky)

VP \rightarrow V NP This is a phrase-structure rule.
 For ex, a verb phrase can be created with just one
 VP \rightarrow V

NP \rightarrow NP Adj \rightarrow phase structure rules based on
 NP \rightarrow N Adj \rightarrow concatenation, and they are
 recursive

They can be repeated as a tree.

There is a symbol denoting the entire input phrase.
 \rightarrow [S \rightarrow NP VP] main structure rule.

These rules can refer both to syntactic rules and
 individual words (by adjectives, lexical symbols).

NP \rightarrow N Adj

Noun Adjective

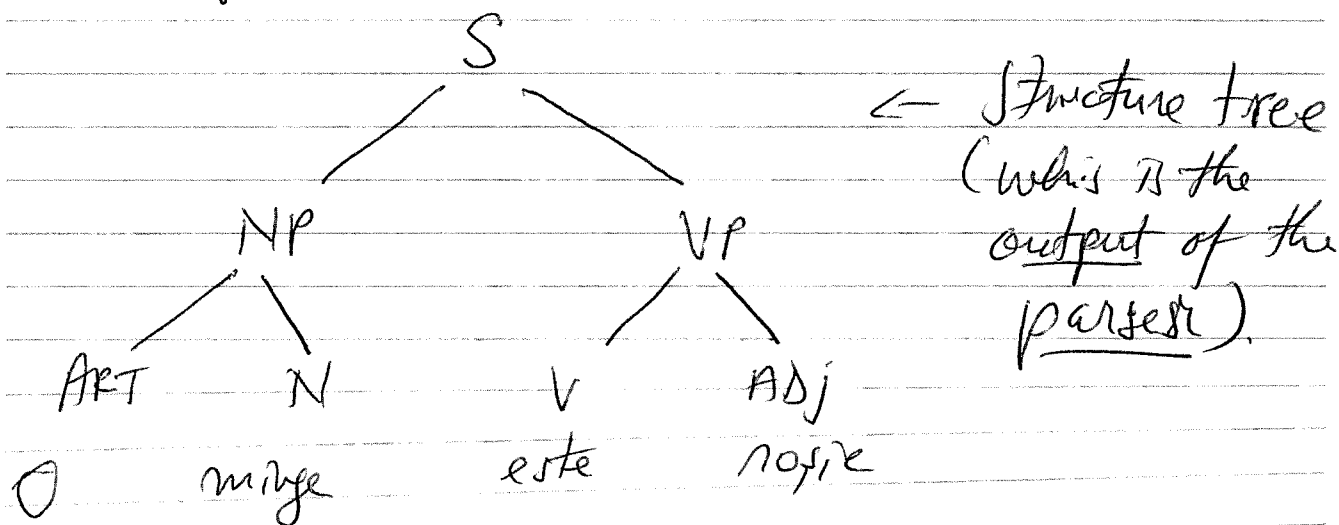
Individual words \rightarrow context words (contribute in
 context)

←
 Nouns
 adjectives
 verbs
 adverbs

Articles, conjunctions, prepositions are words without
 context. Typically, we will insert the leftmost
 to eliminate words without context.

The operation is called parsing, the output of the
 parser is the syntactic structure of the phrase,
 represented by its constituents.

Assuming we have the following inputs:
"O minge este noşie".



$S [np [(art, o), (n, minge)], vp [\dots]]$

↳ list structure that represents the tree.

This represents a syntactical analysis of this phrase.

The main difficulty in natural language is ambiguity.

Ambiguity can be of different types.

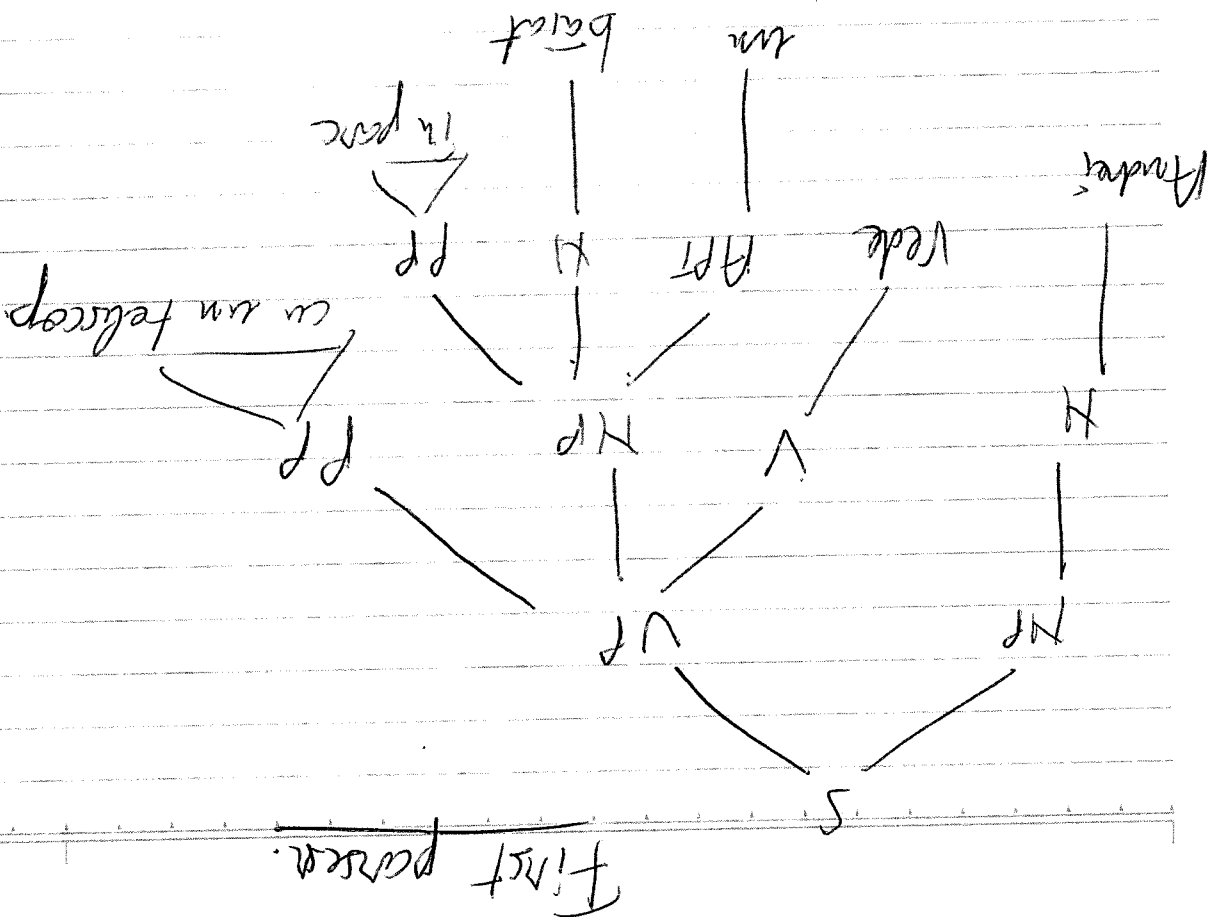
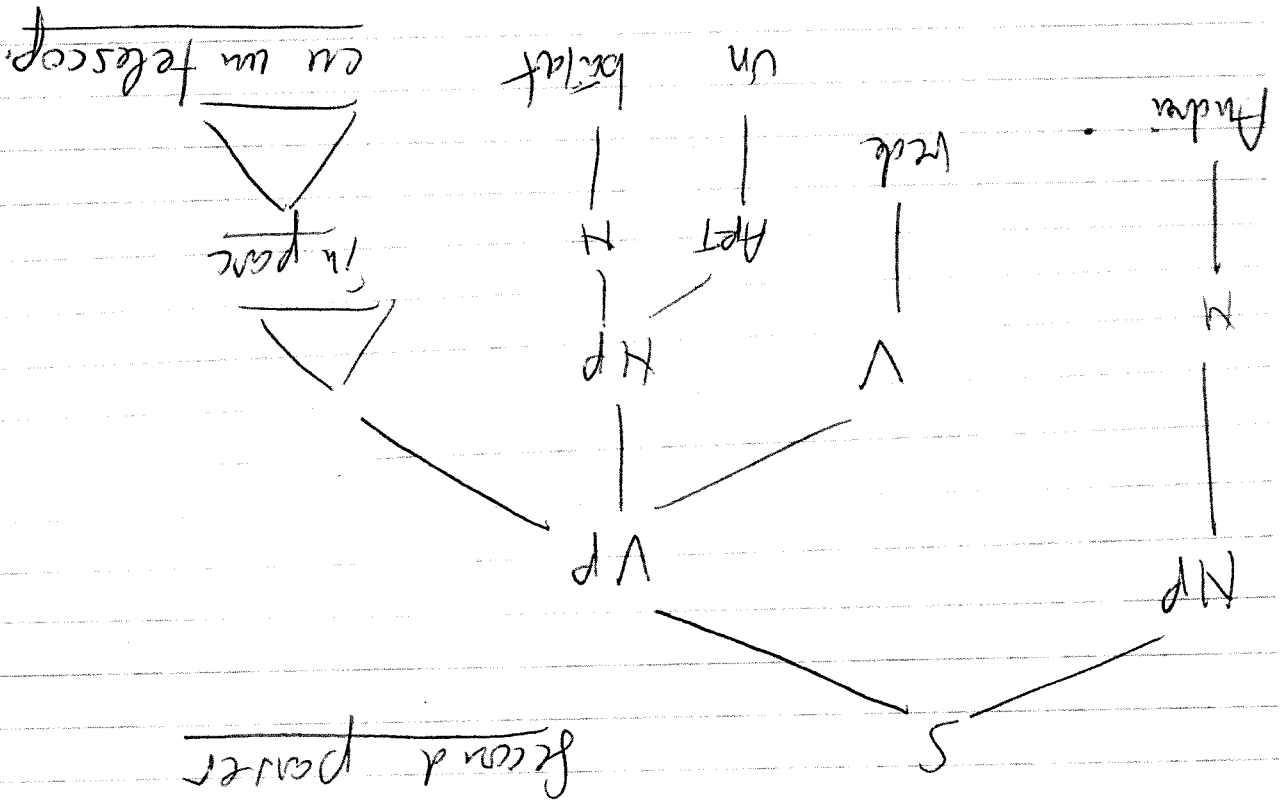
1). Structural ambiguity → This occurs whenever a parser detects different possible structure for the phrase → for each structure there is a different tree. This takes you to the next level.

"Audei vedea un băiat în parcu cu un telescop."

Structural ambiguity leads to semantic ambiguity.

An intelligent algorithm should give me two parsing trees.

2). Lexical ambiguity → a word can be different parts of speech. (ex: nice flies).
 Lexical ambiguity can lead to structural ambiguity.



Sometimes, ambiguity can be resolved.

Ex. Flying planes are dangerous.

Flying planes is dangerous.

↳ the singular gives us the
part of speech Verb.

3. Semantic ambiguity, which concerns us the most.

A word has multiple meanings. A word has several meanings.

Ex of a sentence which is not structurally ambiguous, but it is semantically ambiguous.

"Fiecare copil areste o prajitura"

Natural Language Processing - Laboratory

dir(text 2)

text 2. vocab(). keys()

l.sort()

print(l)

[exp(x) for x in list if cond(x)]

string.isalnum() → The isalnum() method

returns True if all characters in the string are alphanumeric.

If not, returns False.

"abc".isalnum()