

CSCE 771: Computer Processing of Natural Language

Lecture 6: Shallow Parsing, Dependency Parsing, Project Review

PROF. BIPLAV SRIVASTAVA, AI INSTITUTE

5TH SEPTEMBER, 2024

Carolinian Creed: “I will practice personal and academic integrity.”

Acknowledgement: Used materials by
Profs. Mausam, Jurafsky & Martin,
Stanford NLP

Organization of Lecture 6

- Opening Segment
 - Review of Last Lecture
- Main Lecture
- Concluding Segment
 - About Next Lecture – Lecture 7



Main Section

- Project: complete reviewing of topics
- Shallow Parsing
- Dependency Parsing

Recap of Lecture 5

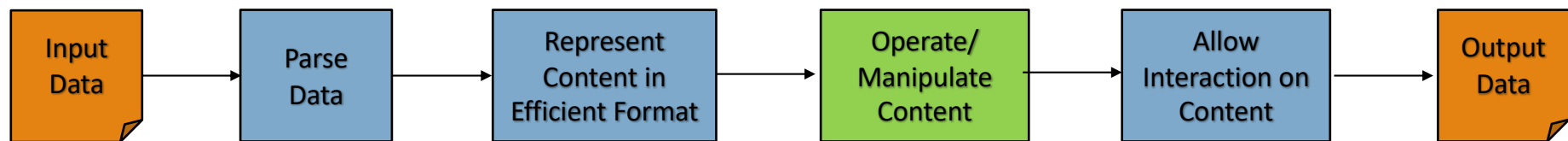
- We discussed the paper - "Contextual Word Representations: Putting Words into Computers", by Noah Smith, CACM June 2020"
- We looked at parsing
 - Roles it plays: verifying , generating, recognizing
 - Many types of parsing: shallow parsing for quick NLP tasks, phrase structure parsing, dependency parsing
- Started reviewing projects

Announcements

- Quiz 1 - next class, in-person and using paper and pen
 - No makeup (best of 3 from 4 quizzes)
 - Will cover concepts discussed in class

Main Lecture

Review: Parsing



Parsing

Types of Parsing

- **Phrase structure / Constituency Parsing:** find phrases and their recursive structure.
Constituency - groups of words behaving as single units, or constituents. Context free grammars are also called Phrase-Structure Grammars
 - **Shallow Parsing/ Chunking:** identify the flat, non-overlapping segments of a sentence: noun phrases, verb phrases, adjective phrases, and prepositional phrases.
- **Dependency Parsing:** find relations in sentences
- **Probabilistic Parsing:** given a sentence X, predict the most **probable** parse tree Y

Chunking

- Chunking - process of identifying and classifying the flat, non-overlapping segments of a sentence that constitute the basic **non-recursive phrases** corresponding to the major content-word parts-of-speech:

- noun phrases
- verb phrases
- adjective phrases, and
- prepositional phrases

Example

[*NP* The morning flight] [*PP* from] [*NP* Denver] [*VP* has arrived.]

- Two operations in this type of parsing:
 - segmenting - finding the non-overlapping extents of the chunks and
 - labeling - assigning the correct tag to the discovered chunks
- Some words may not be part of any chunk

From Jurafsky & Martin

Shallow Parsing/ Chunking

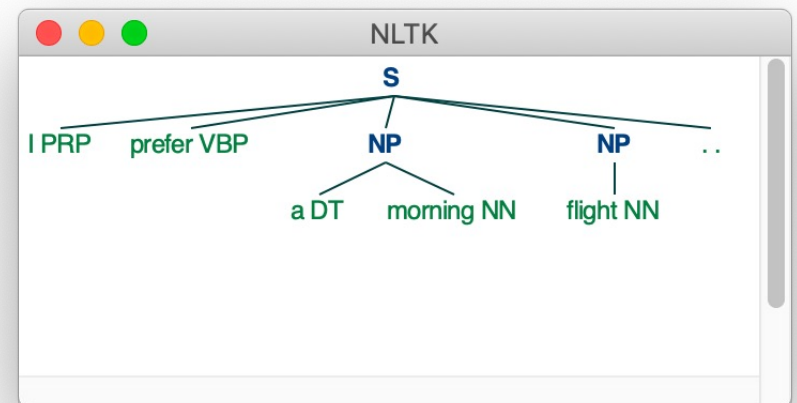
```
data = "I prefer a morning flight."

# Prepare data
tokens = nltk.word_tokenize(data)
tag = nltk.pos_tag(tokens)

# Grammar to use
grammar = "NP: {<DT>?<JJ>*<NN>}"
cp = nltk.RegexpParser(grammar)

# Parse based on regex
result = cp.parse(tag)
print(result)
```

(S I/PRP prefer/VBP (NP a/DT morning/NN) (NP flight/NN) ./.)



IOB notation

- Chunking - IOB tagging
 - B - beginning of each chunk type
 - I - inside of each chunk type
 - O - one for tokens outside (O) any chunk
- Total: $(2N + 1)$ tags for N chunk types

Example

The morning flight from Denver has arrived.
B_NP I_NP I_NP O B_NP O O

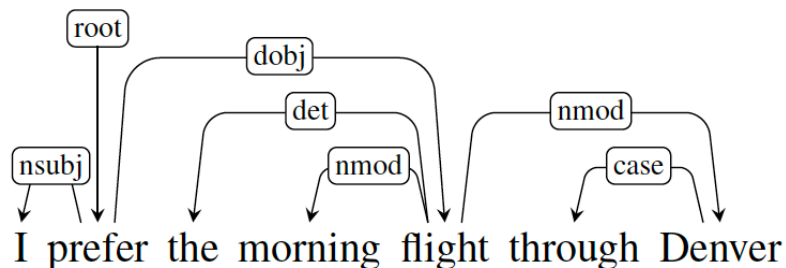
From Jurafsky & Martin

Code and Examples

- Sample code –
<https://github.com/biplav-s/course-nl-f22/blob/main/sample-code/l6-l7-parsing/Chunking%20-%20syntax%20exploration.ipynb>
- Advanced examples –
<https://www.nltk.org/book/ch07.html>

Dependency Parsing

- **Meaning** depends on
 - Words (lemmas) in a sentence
 - Their directed binary grammatical relations with other words
 - (and not on CFGs)
- **Notation:** Labeled arcs are from heads to dependents



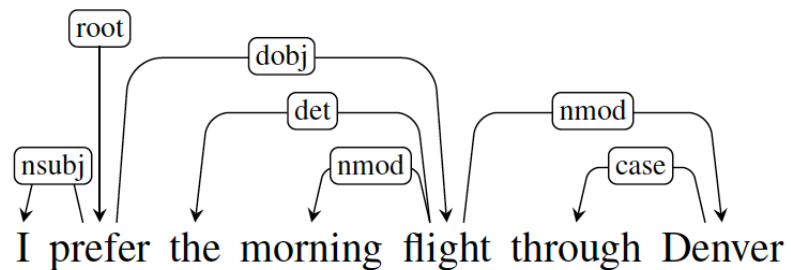
No node corresponding to phrasal constituents or lexical categories in the dependency parse

Dependency Conditions

1. There is a single designated root node that has no incoming arcs.
2. With the exception of the root node, each vertex has exactly one incoming arc.
3. There is a unique path from the root node to each vertex in V.

From Jurafsky & Martin

Dependency Parsing

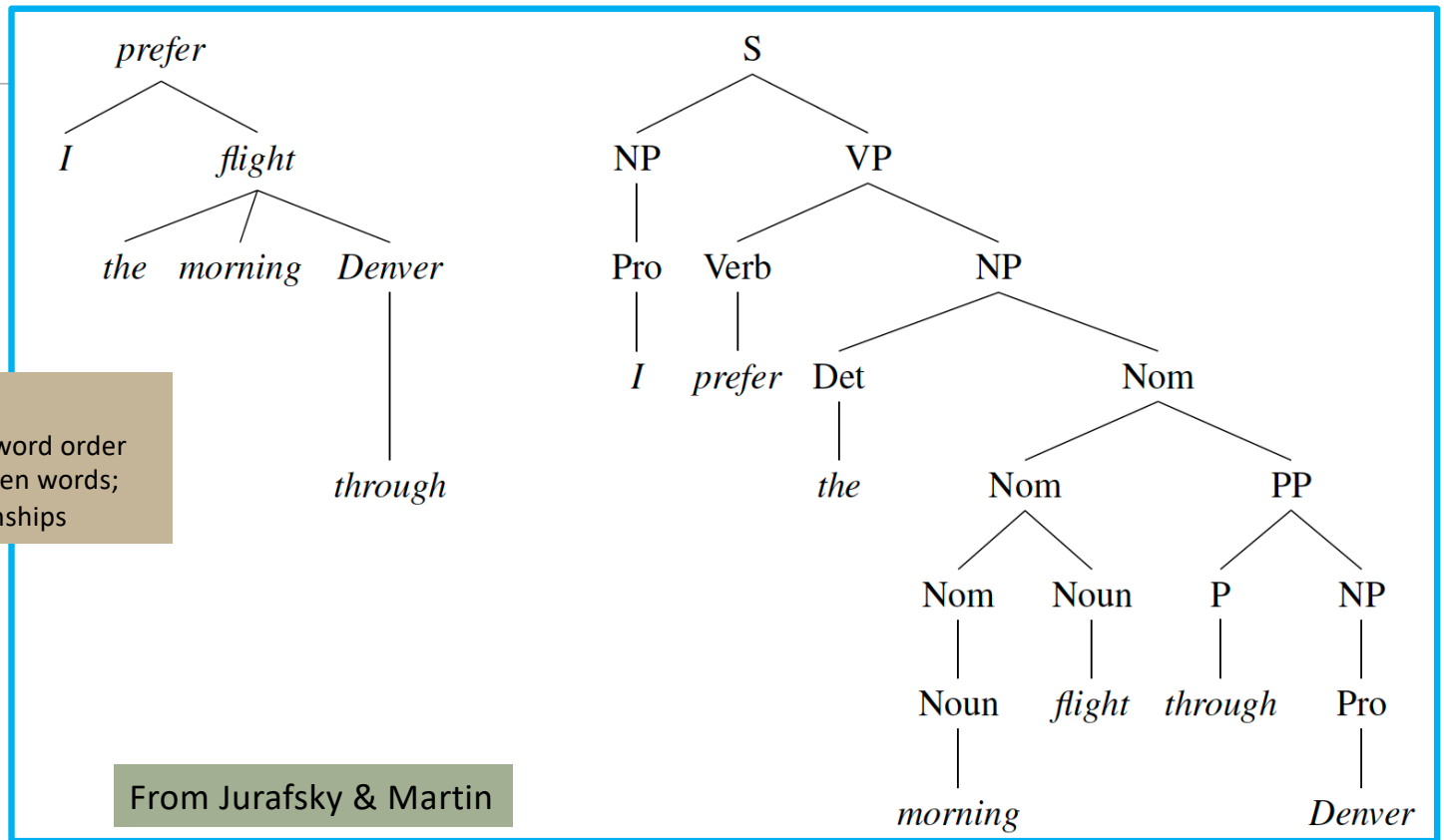


Edge: role that the dependent plays with respect to its head. Examples: subject, direct object and indirect object.

| Clausal Argument Relations | Description |
|----------------------------|--|
| NSUBJ | Nominal subject |
| DOBJ | Direct object |
| IOBJ | Indirect object |
| CCOMP | Clausal complement |
| XCOMP | Open clausal complement |
| Nominal Modifier Relations | Description |
| NMOD | Nominal modifier |
| AMOD | Adjectival modifier |
| NUMMOD | Numeric modifier |
| APPOS | Appositional modifier |
| DET | Determiner |
| CASE | Prepositions, postpositions and other case markers |
| Other Notable Relations | Description |
| CONJ | Conjunct |
| CC | Coordinating conjunction |

From Jurafsky & Martin

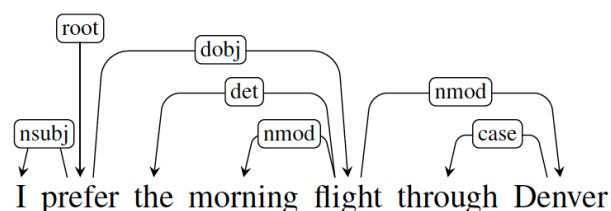
Comparison: Dependency and Phrase Structure



Dependency:

- Useful in languages with free word order
- Highlights relationships between words; useful in tasks needing relationships

Example Dependency Relationships



From Jurafsky & Martin

| Clausal Argument Relations | Description |
|----------------------------|--|
| NSUBJ | Nominal subject |
| DOBJ | Direct object |
| IOBJ | Indirect object |
| CCOMP | Clausal complement |
| XCOMP | Open clausal complement |
| Nominal Modifier Relations | Description |
| NMOD | Nominal modifier |
| AMOD | Adjectival modifier |
| NUMMOD | Numeric modifier |
| APPOS | Appositional modifier |
| DET | Determiner |
| CASE | Prepositions, postpositions and other case markers |
| Other Notable Relations | Description |
| CONJ | Conjunct |
| CC | Coordinating conjunction |

Examples of Parsing with Spacy

- Sample code –
<https://github.com/biplav-s/course-nl-f22/blob/main/sample-code/l6-l7-parsing/parsing%20spacy.ipynb>

See GitHub

About Grammar Forms

- **Strong equivalent grammars:** Two grammars are strongly equivalent if they generate the same set of strings and if they assign the same phrase structure to each sentence
- **Weakly equivalent grammars:** Two grammars are weakly equivalent if they generate the same set of strings but do not assign the same phrase structure to each sentence.
- **Chomsky Normal Form:** a grammar which is
 - ϵ -free and
 - each production is either of the form $A \rightarrow BC$ or $A \rightarrow a$.
- ***Any context-free grammar can be converted into a weakly equivalent Chomsky normal form grammar***
 - $A \rightarrow BCD$ becomes
 - $A \rightarrow BX$ and $X \rightarrow CD$

From Jurafsky & Martin

Review: Parsing - CFG

N a set of **non-terminal symbols** (or **variables**)
 Σ a set of **terminal symbols** (disjoint from N)
 R a set of **rules** or productions, each of the form $A \rightarrow \beta$,
where A is a non-terminal,
 β is a string of symbols from the infinite set of strings $(\Sigma \cup N)^*$
 S a designated **start symbol** and a member of N

Example CFG:

- $N = \{S, NP, VP, \}$
- $\Sigma = \{\text{he, she, walks, sleeps}\}$
- $R = \{$
 - $S \rightarrow NP, VP$
 - $NP \rightarrow \text{he}$
 - $NP \rightarrow \text{she}$
 - $VP \rightarrow \text{walks}$
 - $VP \rightarrow \text{sleeps}$ $\}$
- $S = S$

Questions: *which strings are in the language of example CFG*

(a) she sleeps (b) walks sheeps (c) sleeps he (d) she walks (e) he and she walks

Parsing Perspective

- **Question:** Is parsing of a sentence unique ?

Example 1: “*Book the dinner flight*”

- Book the flight which has dinner
- Book the flight for dinner

Example 2: “*I made her duck*”

- I cooked duck (**sense: animal**) for her
- I cooked duck (**sense: animal**) belonging to her.
- I turned her into duck (**sense: animal**)
- I created duck (**sense: object**) for her
- I made her to lower her head or body (**sense: posture**).

From Jurafsky & Martin

Parsing Perspective

- **Question:** Is parsing of a sentence unique ?
- **Answer:** Not necessarily

Issue: Then, which one to return?

Solution: Given a sentence X, predict its parse tree Y

Lecture 6: Concluding Comments

- We reviewed projects
- We reviewed parsers
 - Shallow parsers
 - Dependency parsers

Concluding Segment

Discussion: Course Project

Theme: Analyze quality of official information available for elections in 2024 [in a state]

- Take information available from
 - Official site: State Election Commissions
 - Respected non-profits: League of Women Voters
- Analyze information
 - State-level: Analyze quality of questions, answers, answers-to-questions
 - Comparatively: above along all states (being done by students)
- Benchmark and report
 - Compare analysis with LLM
 - Prepare report

- Process and analyze using NLP
 - Extract entities
 - Assess quality – metrics
 - Content – *Englishness*
 - Content – *Domain* -- election
 - ... other NLP tasks
- Analyze and communicate overall

Major dates for project check

- Sep 10: written – project outline
- Oct 8: in class
- Oct 31: in class // LLM
- Dec 5: in class // Comparative

Review current states chosen by others

Project Discussion

1. Go to Google spreadsheet against your name
2. Enter the state you will focus on for course project

1. Create a private Github repository called “CSCE771-Fall2024-<studentname>-Repo”. Share with Instructor (biplav-s) and TA (vr25)
2. Create Google folder called “CSCE771-Fall2024-<studentname>-SharedInfo”. Share with Instructor (prof.biplav@gmail.com) and TA (rawtevipula25@gmail.com)
3. Create a Google doc in your Google repo called “Project Plan” and have the following by Friday (Aug 30, 2024)

Timeline

1. Title: [Analyze quality of official information available for elections in 2024](#) in <state>
2. Data need:
 1. Official: state’s election commission
 2. LWV: <https://www.vote411.org/>
3. Methods:
4. Evaluation:
5. Milestones
 - Sep 10: written and feedback
 - Oct 8: in class
 - Oct 31: in class
 - Dec 5: in class

Discussion: Course Project

- **Expectations**

- Apply methods learned in class or of interest to a problem of interest
- Be goal oriented: aim to finish, be proactive, be innovative
- Do top-class work: code, writeup, presentation

- **Typical pitfalls**

- Not detailing out the project, assuming data
- Not spending enough time

- **What will be awarded**

- Results and efforts (balance)
- Challenge level of problem

Review current states chosen by others

Course Project – Deadlines and Penalty Rubric

- Penalty
 - Missing milestones: [-10%]
 - Maximum: [-40%]
- Bonus possible
 - if two or more states considered
 -

Timeline

1. Title: Analyze quality of official information available for elections in 2024 in <state>
2. Data need:
 1. Official: state's election commission
 2. LWV: <https://www.vote411.org/>
3. Methods:
4. Evaluation:
5. Milestones
 - Sep 10: written and feedback
 - Oct 8: in class
 - Oct 31: in class
 - Dec 5: in class

About Next Lecture – Lecture 7

Lecture 7

- Statistical parsing
- QUIZ

| | | | | |
|----|-------------|---|--|--|
| 4 | Aug 29 (Th) | NLP Tasks, Case Study – Business Application | | |
| 5 | Sep 3 (Tu) | Parsing, Paper 1 discussion; project topics review | | Practice exercise |
| 6 | Sep 5 (Th) | Project topics review, statistic Parsing | | |
| 7 | Sep 10 (Tu) | Statistical parsing, QUIZ | | Quiz 1, Project Check |
| 8 | Sep 12 (Th) | Evaluation, Semantics | | Coding running example |
| 9 | Sep 17 (Tu) | Semantics Machine Learning for NLP, Evaluation - Metrics | | Code: scikit fl score package, Code: ConceptIO |
| 10 | Sep 19 (Th) | Towards Language Model: Vector embeddings, Embeddings, CNN/ RNN | | Code: embedding, genism word vector, tf-idf |