# Speech and Language Processing

Chapter 12
Syntactic Parsing

# Today

- Parsing with CFGs
  - Bottom-up, top-down
  - Ambiguity
  - CKY parsing
  - (Earley)
  - Shallow

# Parsing

- Parsing with CFGs refers to the task of assigning proper trees to input strings
- Proper here means a tree that covers all and only the elements of the input and has an S at the top
- It doesn't actually mean that the system can select the correct tree from among all the possible trees

01/29/20      Speech and Language Processing - Jurafsky and Martin      3

# Parsing

- As with everything of interest, parsing involves a search which involves the making of choices
- We'll start with some basic (meaning bad) methods before moving on to the one that you need to know

01/29/20      Speech and Language Processing - Jurafsky and Martin      4

# For Now

- Assume…
  - You have all the words already in some buffer
  - The input isn't POS tagged
  - We won't worry about morphological analysis
  - All the words are known

  - These are all problematic in various ways, and would have to be addressed in real applications.

# Top-Down Search

- Since we're trying to find trees rooted with an *S* (Sentences), why not start with the rules that give us an *S*.
- Then we can work our way down from there to the words.

- "Book that flight"

# Bottom-Up Parsing

- Of course, we also want trees that cover the input words. So we might also start with trees that link up with the words in the right way.
- Then work your way up from there to larger and larger trees.

Speech and Language Processing - Jurafsky and Martin

# *"The old dog the footsteps of the young."*

| | |
|---|---|
| S → NP VP | VP → V |
| S → Aux NP VP | VP -> V PP |
| S -> VP | PP -> Prep NP |
| NP → Det Nom | N → old \| dog \| footsteps \| young |
| NP →PropN | V → dog \| eat \| sleep \| bark \| meow |
| Nom -> Adj N | Aux → does \| can |
| Nom → N | Prep →from \| to \| on \| of |
| Nom → N Nom | PropN → Fido \| Felix |
| Nom → Nom PP | Det → that \| this \| a \| the |
| VP → V NP | Adj -> old \| happy\| young |

# Top-Down and Bottom-Up

- Top-down
  - Only searches for trees that can be answers (i.e. S's)
  - But also suggests trees that are not consistent with any of the words
- Bottom-up
  - Only forms trees consistent with the words
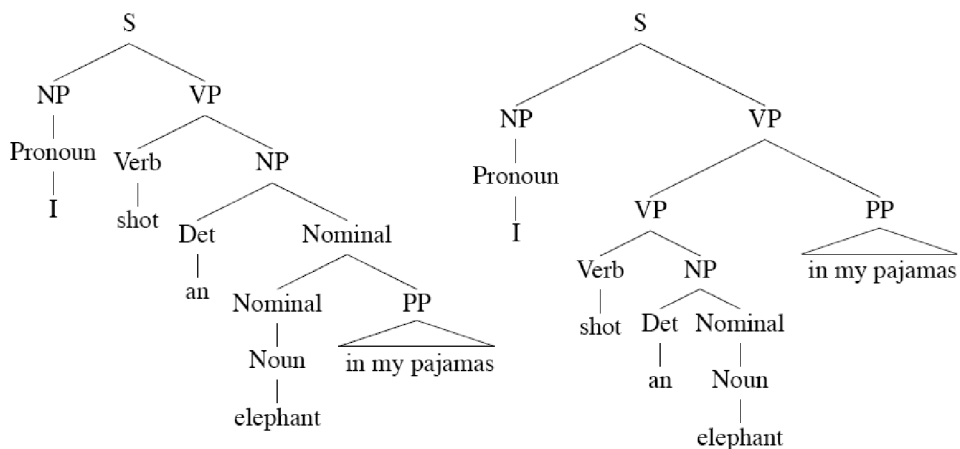  - But suggests trees that make no sense globally

# Control

- Of course, in both cases we left out how to keep track of the search space and how to make choices
  - Which node to try to expand next
  - Which grammar rule to use to expand a node
- One approach is called backtracking.
  - Make a choice, if it works out then fine
  - If not then back up and make a different choice

# Problems

- Even with the best filtering, backtracking methods are doomed because of two inter-related problems
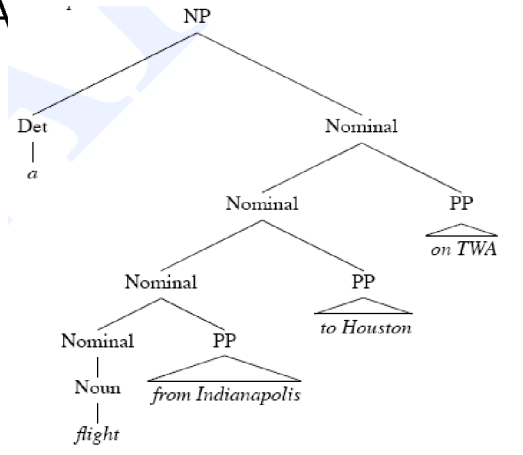  - Ambiguity
  - Shared subproblems

Speech and Language Processing - Jurafsky and Martin

# Ambiguity



Speech and Language Processing - Jurafsky and Martin

# Example types of ambiguity

- POS
- Attachment
  - PP
  - Coordination (*old dogs and cats*)

Speech and Language Processing - Jurafsky and Martin

# Shared Sub-Problems

- No matter what kind of search (top-down or bottom-up or mixed) that we choose.
  - We don't want to redo work we've already done.
  - Unfortunately, naïve backtracking will lead to duplicated work.

Speech and Language Processing - Jurafsky and Martin

# Shared Sub-Problems

- Consider
  - A flight from Indianapolis to Houston on TWA

# Shared Sub-Problems

- Assume a top-down parse making choices among the various Nominal rules.
- In particular, between these two
  - Nominal -> Noun
  - Nominal -> Nominal PP
- Statically choosing the rules in this order leads to the following bad results...
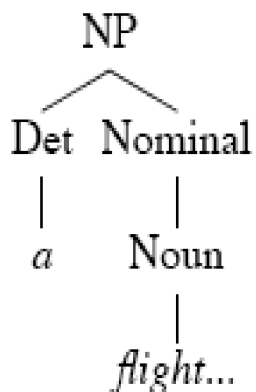
# Shared Sub-Problems

```
              NP
             /  \
          Det    Nominal
           |        |
           a      Noun
                    |
                  flight...
```

Speech and Language Processing - Jurafsky and Martin

# Shared Sub-Problems

```
                NP
               /  \
            Det    Nominal
             |      /    \
             a   Nominal   PP
                    |      /  \
                  Noun   from Indianapolis...
                    |
                 flight
```
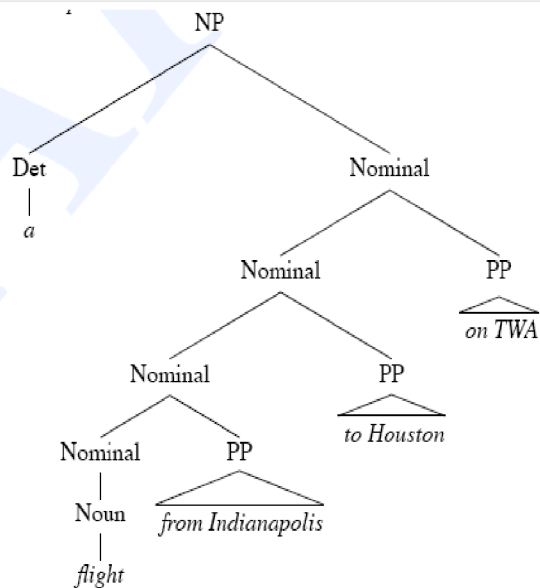
Speech and Language Processing - Jurafsky and Martin

# Shared Sub-Problems



Speech and Language Processing - Jurafsky and Martin

# Shared Sub-Problems



Speech and Language Processing - Jurafsky and Martin

# Dynamic Programming

- DP search methods fill tables with partial results and thereby
  - Avoid doing avoidable repeated work
  - Solve exponential problems in polynomial time
  - Efficiently store ambiguous structures with shared sub-parts.
- Two approaches roughly correspond to bottom-up and top-down approaches.
  - CKY
  - Earley

01/29/20          Speech and Language Processing - Jurafsky and Martin          27

# CKY Parsing

- First we'll limit our grammar to epsilon-free, binary rules (more later)

- Consider the rule $A \rightarrow BC$
  - If there is an A somewhere in the input then there must be a B followed by a C in the input.
  - If the A spans from i to j in the input then there must be some k st. i<k<j
    - Ie. The B splits from the C someplace.

01/29/20          Speech and Language Processing - Jurafsky and Martin          28

# Problem

- What if your grammar isn't binary?
  - As in the case of the TreeBank grammar?
- Convert it to binary… any arbitrary CFG can be rewritten into Chomsky-Normal Form automatically.
- What does this mean?
  - The resulting grammar accepts (and rejects) the same set of strings as the original grammar.
  - But the resulting derivations (trees) are different.

01/29/20        Speech and Language Processing - Jurafsky and Martin        29

# Problem

- More specifically, we want our rules to be of the form

  A → B C

  Or

  A → w

  *That is, rules can expand to either 2 non-terminals or to a single terminal.*

01/29/20        Speech and Language Processing - Jurafsky and Martin        30

# Binarization Intuition

- Eliminate chains of unit productions.
- Introduce new intermediate non-terminals into the grammar that distribute rules with length > 2 over several rules.
  - So... *S → A B C turns into*

  *S → X C and*

  *X → A B*

  Where X is a symbol that doesn't occur anywhere else in the the grammar.

# Sample L1 Grammar

| Grammar | Lexicon |
|---|---|
| $S \rightarrow NP\ VP$ | $Det \rightarrow that \mid this \mid a$ |
| $S \rightarrow Aux\ NP\ VP$ | $Noun \rightarrow book \mid flight \mid meal \mid money$ |
| $S \rightarrow VP$ | $Verb \rightarrow book \mid include \mid prefer$ |
| $NP \rightarrow Pronoun$ | $Pronoun \rightarrow I \mid she \mid me$ |
| $NP \rightarrow Proper\text{-}Noun$ | $Proper\text{-}Noun \rightarrow Houston \mid NWA$ |
| $NP \rightarrow Det\ Nominal$ | $Aux \rightarrow does$ |
| $Nominal \rightarrow Noun$ | $Preposition \rightarrow from \mid to \mid on \mid near \mid through$ |
| $Nominal \rightarrow Nominal\ Noun$ | |
| $Nominal \rightarrow Nominal\ PP$ | |
| $VP \rightarrow Verb$ | |
| $VP \rightarrow Verb\ NP$ | |
| $VP \rightarrow Verb\ NP\ PP$ | |
| $VP \rightarrow Verb\ PP$ | |
| $VP \rightarrow VP\ PP$ | |
| $PP \rightarrow Preposition\ NP$ | |

# CNF Conversion

| $\mathscr{L}_1$ Grammar | $\mathscr{L}_1$ in CNF |
|---|---|
| $S \rightarrow NP\ VP$ | $S \rightarrow NP\ VP$ |
| $S \rightarrow Aux\ NP\ VP$ | $S \rightarrow X1\ VP$ |
| | $X1 \rightarrow Aux\ NP$ |
| $S \rightarrow VP$ | $S \rightarrow book \mid include \mid prefer$ |
| | $S \rightarrow Verb\ NP$ |
| | $S \rightarrow X2\ PP$ |
| | $S \rightarrow Verb\ PP$ |
| | $S \rightarrow VP\ PP$ |
| $NP \rightarrow Pronoun$ | $NP \rightarrow I \mid she \mid me$ |
| $NP \rightarrow Proper\text{-}Noun$ | $NP \rightarrow TWA \mid Houston$ |
| $NP \rightarrow Det\ Nominal$ | $NP \rightarrow Det\ Nominal$ |
| $Nominal \rightarrow Noun$ | $Nominal \rightarrow book \mid flight \mid meal \mid money$ |
| $Nominal \rightarrow Nominal\ Noun$ | $Nominal \rightarrow Nominal\ Noun$ |
| $Nominal \rightarrow Nominal\ PP$ | $Nominal \rightarrow Nominal\ PP$ |
| $VP \rightarrow Verb$ | $VP \rightarrow book \mid include \mid prefer$ |
| $VP \rightarrow Verb\ NP$ | $VP \rightarrow Verb\ NP$ |
| $VP \rightarrow Verb\ NP\ PP$ | $VP \rightarrow X2\ PP$ |
| | $X2 \rightarrow Verb\ NP$ |
| $VP \rightarrow Verb\ PP$ | $VP \rightarrow Verb\ PP$ |
| $VP \rightarrow VP\ PP$ | $VP \rightarrow VP\ PP$ |
| $PP \rightarrow Preposition\ NP$ | $PP \rightarrow Preposition\ NP$ |

# CKY

- So let's build a table so that an A spanning from i to j in the input is placed in cell [i,j] in the table.
- So a non-terminal spanning an entire string will sit in cell [0, n]
  - Hopefully an S
- If we build the table bottom-up, we'll know that the parts of the A must go from i to k and from k to j, for some k.