

Natural Language Processing

Lecture 15—10/15/2015
Jim Martin

Today

- Start on Parsing
 - Parsing frameworks
 - CKY

10/15/15

Speech and Language Processing - Jurafsky and Martin

2

Treebanks

- Treebanks are corpora in which each sentence has been paired with a parse tree (presumably the right one).
- These are generally created
 1. By first parsing the collection with an automatic parser
 2. And then having human annotators hand correct each parse as necessary.
- This generally requires detailed annotation guidelines that provide a POS tagset, a grammar, and instructions for how to deal with particular grammatical constructions.

10/15/15

Speech and Language Processing - Jurafsky and Martin

3

Head Finding

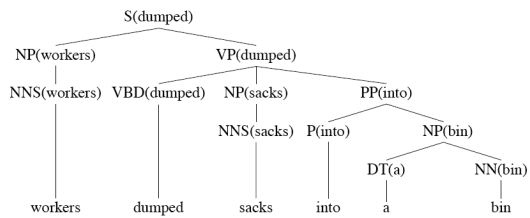
- Finding heads in treebank trees is a task that arises frequently in many applications.
 - As we'll see it is particularly important in statistical parsing
- We can visualize this task by annotating the nodes of a parse tree with the heads of each corresponding node.

10/15/15

Speech and Language Processing - Jurafsky and Martin

7

Lexically Decorated Tree



10/15/15

Speech and Language Processing - Jurafsky and Martin

8

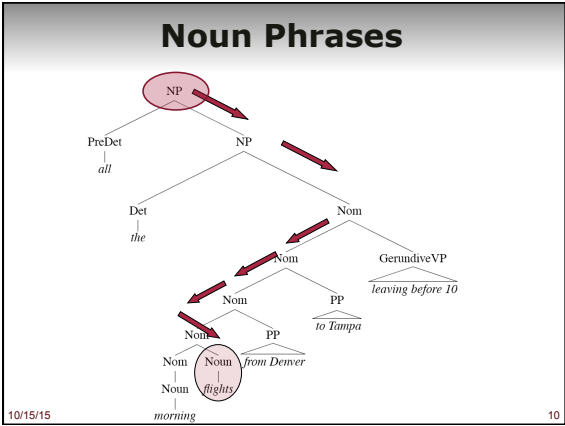
Head Finding

- Given a tree, the standard way to do head finding is to use a simple set of tree traversal rules specific to each non-terminal in the grammar.

10/15/15

Speech and Language Processing - Jurafsky and Martin

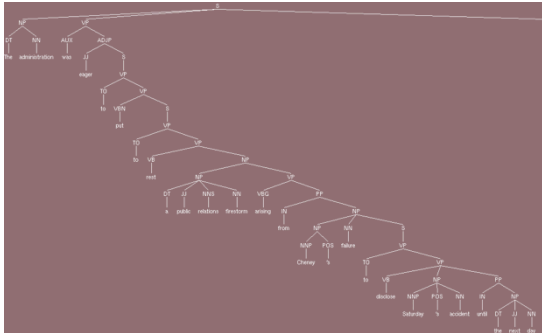
9



- ### Treebank Uses
- Treebanks (and head-finding) are particularly critical to the development of statistical parsers
 - Chapter 14
 - Also valuable to *Corpus Linguistics*
 - Investigating the empirical details of various constructions in a given language
- 10/15/15 11

- ### 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 mean that the system can select the correct tree from among all the possible trees
- 10/15/15 12

Automatic Syntactic Parse



For Now

- Let's assume...
 - You have all the words for a sentence already in some buffer
 - The input is not POS tagged prior to parsing
 - 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.

10/15/15

Speech and Language Processing - Jurafsky and Martin

14

Search Framework

- It's productive to think about parsing as a form of search...
 - A search through the space of possible trees given an input sentence and grammar
 - This framework suggests that heuristic search methods and/or dynamic programming methods might be applicable
 - It also suggests that notions such as the direction of the search might be useful

10/15/15

Speech and Language Processing - Jurafsky and Martin

15

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.

10/15/15 Speech and Language Processing - Jurafsky and Martin 16

Top Down Space

10/15/15 Speech and Language Processing - Jurafsky and Martin 17

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.

10/15/15 Speech and Language Processing - Jurafsky and Martin 18

Bottom-Up Search

Book that flight

10/15/15 Speech and Language Processing - Jurafsky and Martin 19

Bottom-Up Search

Verb Det Noun
| | |
Book that flight

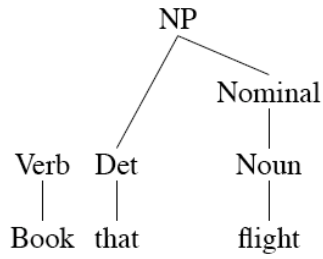
10/15/15 Speech and Language Processing - Jurafsky and Martin 20

Bottom-Up Search

Nominal
|
Noun
| | |
Verb Det Noun
| | |
Book that flight

10/15/15 Speech and Language Processing - Jurafsky and Martin 21

Bottom-Up Search

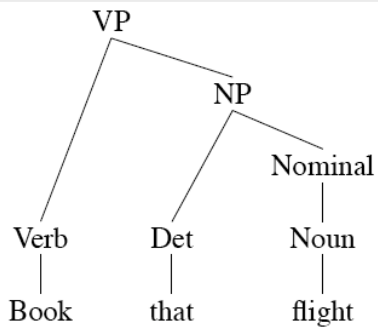


10/15/15

Speech and Language Processing - Jurafsky and Martin

22

Bottom-Up Search



10/15/15

Speech and Language Processing - Jurafsky and Martin

23

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

10/15/15

Speech and Language Processing - Jurafsky and Martin

24

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
 - Same as with ND-Recognize

10/15/15

Speech and Language Processing - Jurafsky and Martin

25

Problems

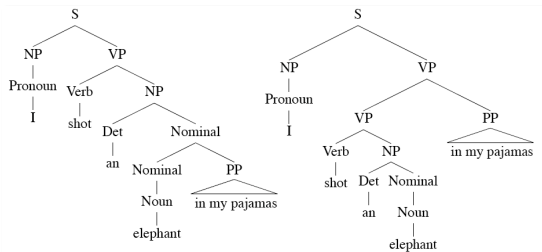
- Even with the best filtering, backtracking methods are doomed because of two inter-related problems
 - Ambiguity and search control (choice)
 - Shared subproblems

10/15/15

Speech and Language Processing - Jurafsky and Martin

26

Ambiguity



10/15/15

Speech and Language Processing - Jurafsky and Martin

27

Shared Sub-Problems

- No matter what kind of search (top-down or bottom-up or mixed) that we choose...
 - We can't afford to redo work we've already done.
 - Without some help naïve backtracking will lead to such duplicated work.

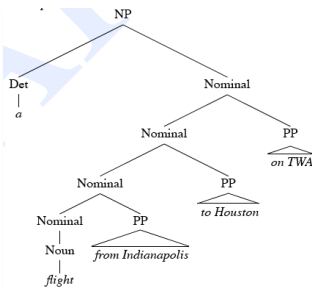
10/15/15

Speech and Language Processing - Jurafsky and Martin

28

Shared Sub-Problems

- Consider
 - *A flight from Indianapolis to Houston on TWA*



10/15/15

Speech and Language Processing - Jurafsky and Martin

29

Sample L1 Grammar

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

10/15/15

Speech and Language Processing - Jurafsky and Martin

30

Shared Sub-Problems

- Assume a top-down parse that has already expanded the *NP* rule (dealing with the Det)
- Now its 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 behavior...

10/15/15 Speech and Language Processing - Jurafsky and Martin 31

Shared Sub-Problems

10/15/15 Speech and Language Processing - Jurafsky and Martin 32

Shared Sub-Problems

10/15/15 Speech and Language Processing - Jurafsky and Martin 33

Shared Sub-Problems

The diagram shows a parse tree for the sentence "a flight from Indianapolis to Houston...". The root node is NP, which branches into Det (a) and Nominal. This Nominal node branches into another Nominal and PP (to Houston...). The second Nominal node branches into Noun (flight) and PP (from Indianapolis). The third PP node branches into to Houston... and another PP (from Indianapolis). A red circle with a diagonal slash is drawn over the root NP node and its children, indicating a shared sub-problem.

10/15/15 34
Speech and Language Processing - Jurafsky and Martin

Shared Sub-Problems

The diagram shows a parse tree for the sentence "a flight from Indianapolis to Houston on TWA". The root node is NP, which branches into Det (a) and Nominal. This Nominal node branches into another Nominal and PP (on TWA). The second Nominal node branches into another Nominal and PP (to Houston). The third Nominal node branches into Noun (flight) and PP (from Indianapolis). The fourth PP node branches into to Houston and another PP (on TWA). A red circle with a diagonal slash is drawn over the root NP node and its children, indicating a shared sub-problem.

10/15/15 35
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 (ok, not really)
 - Efficiently store ambiguous structures with shared sub-parts.
- We'll cover one approach that corresponds to a bottom-up strategy
 - CKY

10/15/15 36
Speech and Language Processing - Jurafsky and Martin

CKY Parsing

- First we'll limit our grammar to epsilon-free, binary rules (more on this later)
- Consider the rule $A \rightarrow BC$
 - If there is an A somewhere in the input generated by this rule 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. $i < k < j$
 - In other words, the B splits from the C someplace after the i and before the j .

10/15/15

Speech and Language Processing - Jurafsky and Martin

37

CKY

- 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 it will be an S
- Now we know that the parts of the A must go from i to k and from k to j , for some k

10/15/15

Speech and Language Processing - Jurafsky and Martin

38

CKY

- Meaning that for a rule like $A \rightarrow BC$ we should look for a B in $[i,k]$ and a C in $[k,j]$.
- In other words, if we think there might be an A spanning i,j in the input... AND $A \rightarrow BC$ is a rule in the grammar THEN
- There must be a B in $[i,k]$ and a C in $[k,j]$ for some k such that $i < k < j$

What about the B and the C?

10/15/15

Speech and Language Processing - Jurafsky and Martin

39

CKY

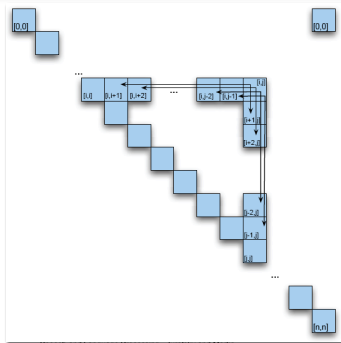
- So to fill the table loop over the cells $[i,j]$ values in some systematic way
 - Then for each cell, loop over the appropriate k values to search for things to add.
 - Add all the derivations that are possible for each $[i,j]$ for each k

10/15/15

Speech and Language Processing - Jurafsky and Martin

40

CKY Table



10/15/15

Speech and Language Processing - Jurafsky and Martin

41

CKY Algorithm

function CKY-PARSE(*words*, *grammar*) **returns** *table*

```
for  $j \leftarrow$  from 1 to LENGTH(words) do  
   $table[j-1, j] \leftarrow \{A \mid A \rightarrow words[j] \in grammar\}$   
  for  $i \leftarrow$  from  $j-2$  downto 0 do  
    for  $k \leftarrow i+1$  to  $j-1$  do  
       $table[i, j] \leftarrow table[i, j] \cup$   
         $\{A \mid A \rightarrow BC \in grammar,$   
           $B \in table[i, k],$   
           $C \in table[k, j]\}$ 
```

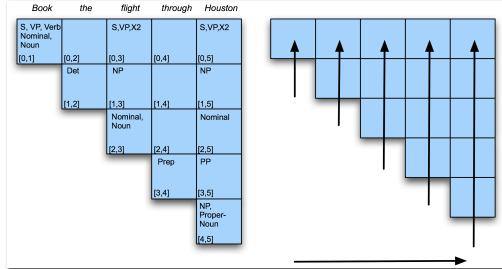
What's the complexity of this?

10/15/15

Speech and Language Processing - Jurafsky and Martin

42

Example

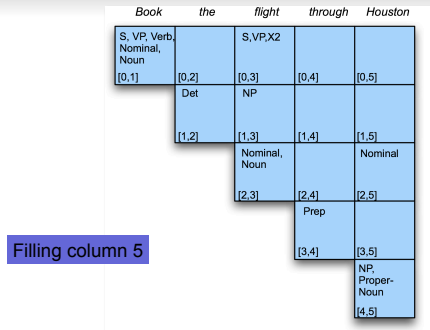


10/15/15

Speech and Language Processing - Jurafsky and Martin

43

Example



Filling column 5

10/15/15

Speech and Language Processing - Jurafsky and Martin

44

Example

- Filling column 5 corresponds to processing word 5, which is *Houston*.
 - So j is 5.
 - So i goes from 3 to 0 (3,2,1,0)

```
function CKY-PARSE(words, grammar) returns table
for j ← from 1 to LENGTH(words) do
  table[j-1, j] ← {A | A → words[j] ∈ grammar}
for i ← from j-2 downto 0 do
  for k ← i+1 to j-1 do
    table[i, j] ← table[i, j] ∪
      {A | A → BC ∈ grammar,
        B ∈ table[i, k],
        C ∈ table[k, j]}
```

10/15/15

Speech and Language Processing - Jurafsky and Martin

45

Example

Book	the	flight	through	Houston
S, VP, Verb, Nominal, Noun [0,1]	[0,2]	S,VP,X2 [0,3]	[0,4]	[0,5]
Det		NP		NP
[1,2]	[1,3]	[1,4]	[1,5]	
	Nominal, Noun			
	[2,3]	[2,4]	[2,5]	
		Prep	PP	
		[3,4]	[3,5]	
				NP, Proper- Noun
				[4,5]

10/15/15

Speech and Language Processing - Jurafsky and Martin

46

Example

Book	the	flight	through	Houston
S, VP, Verb, Nominal, Noun [0,1]	[0,2]	S,VP,X2 [0,3]	[0,4]	[0,5]
Det		NP		NP
[1,2]	[1,3]	[1,4]	[1,5]	
	Nominal, Noun			
	[2,3]	[2,4]	[2,5]	
		Prep	PP	
		[3,4]	[3,5]	
				NP, Proper- Noun
				[4,5]

10/15/15

Speech and Language Processing - Jurafsky and Martin

47

Example

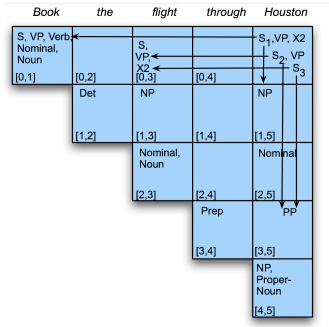
Book	the	flight	through	Houston
S, VP, Verb, Nominal, Noun [0,1]	[0,2]	S,VP,X2 [0,3]	[0,4]	[0,5]
Det		NP		NP
[1,2]	[1,3]	[1,4]	[1,5]	
	Nominal, Noun			
	[2,3]	[2,4]	[2,5]	
		Prep	PP	
		[3,4]	[3,5]	
				NP, Proper- Noun
				[4,5]

10/15/15

Speech and Language Processing - Jurafsky and Martin

48

Example



10/15/15

Speech and Language Processing - Jurafsky and Martin

49

Example

- Since there's an *S* in [0,5] we have a valid parse.
- Are we done? We we sort of left something out of the algorithm

```
function CKY-PARSE(words, grammar) returns table
  for j ← from 1 to LENGTH(words) do
    table[j-1, j] ← {A | A → words[j] ∈ grammar}
  for i ← from j-2 downto 0 do
    for k ← i+1 to j-1 do
      table[i, j] ← table[i, j] ∪
        {A | A → BC ∈ grammar,
          B ∈ table[i, k],
          C ∈ table[k, j]}
```

10/15/15

Speech and Language Processing - Jurafsky and Martin

50

CKY Notes

- Since it's bottom up, CKY hallucinates a lot of silly constituents.
 - Segments that by themselves are constituents but cannot really occur in the context in which they are being suggested.
 - To avoid this we can switch to a top-down control strategy
 - Or we can add some kind of filtering that blocks constituents where they can not happen in a final analysis.

10/15/15

Speech and Language Processing - Jurafsky and Martin

51

CKY Notes

- We arranged the loops to fill the table a column at a time, from left to right, bottom to top.
 - This assures us that whenever we're filling a cell, the parts needed to fill it are already in the table (to the left and below)
 - It's somewhat natural in that it processes the input a left to right a word at a time
 - Known as online
 - Can you think of an alternative strategy?

10/15/15

Speech and Language Processing - Jurafsky and Martin

52
