



IT Systems Engineering | Universität Potsdam

Natural Language Processing

Parsing

Potsdam, 10 May 2012

Saeedeh Momtazi
Information Systems Group

based on the slides of the course book

Parsing

2

- Finding structural relationship between words in a sentence

- Applications
 - Spell checking
 - Speech recognition
 - Machine translation
 - Language modeling

Outline

3

- 1 Phrase Structure
- 2 Syntactic Parsing
CKY Algorithm
- 3 Statistical Parsing

Outline

4

- 1 Phrase Structure
- 2 Syntactic Parsing
CKY Algorithm
- 3 Statistical Parsing

Constituency

5

- Working based on Constituency (Phrase structure)
 - Organizing words into nested constituents
 - Showing that groups of words within utterances can act as single units
 - Forming coherent classes from these units that can behave in similar ways
 - With respect to their internal structure
 - With respect to other units in the language
 - Considering a **head** word for each constituent

Constituency

6

the writer talked to the audiences about his new book.

the writer talked about his new book to the audiences. ✓

about his new book the writer talked to the audiences. ✓

the writer talked book to the audiences about his new. ✗

Context Free Grammar (CFG)

- Grammar G consists of
 - Terminals (T)
 - Non-terminals (N)
 - Start symbol (S)
 - Rules (R)

- **Terminals**
 - The set of words in the text

- **Non-Terminals**
 - The constituents in a language (noun phrase, verb phrase,)

- **Start symbol**
 - The main constituent of the language (sentence)

- **Rules**
 - Equations that consist of a single non-terminal on the left and any number of terminals and non-terminals on the right

CFG

9

 $S \rightarrow NP VP$ $S \rightarrow VP$ $NP \rightarrow N$ $NP \rightarrow Det N$ $NP \rightarrow NP NP$ $NP \rightarrow NP PP$ $VP \rightarrow V$ $VP \rightarrow VP PP$ $VP \rightarrow VP NP$ $PP \rightarrow Prep NP$ $N \rightarrow \text{book}$ $V \rightarrow \text{book}$ $Det \rightarrow \text{the}$ $N \rightarrow \text{flight}$ $Prep \rightarrow \text{through}$ $N \rightarrow \text{Houston}$

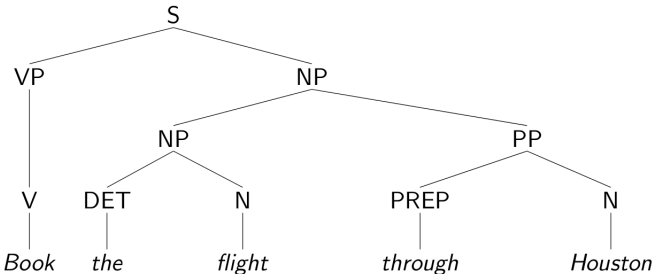
CFG

10

V DET N PREP N
| | | | |
Book *the* *flight* *through* *Houston*

CFG

11



Outline

12

- ① Phrase Structure
- ② Syntactic Parsing
CKY Algorithm
- ③ Statistical Parsing

Parsing

- Taking a string and a grammar and returning proper parse tree(s) for that string
- Covering all and only the elements of the input string
- Reaching the start symbol at the top of the string

- The system cannot select the correct tree among all the possible trees

Main Grammar Fragments

- Sentence
- Noun Phrase
 - Agreement
- Verb Phrase
 - Sub-categorization

Grammar Fragments: Sentence

15

- Declaratives
A plane left.
 $S \rightarrow NP VP$

- Imperatives
Leave!
 $S \rightarrow VP$

- Yes-No Questions
Did the plane leave?
 $S \rightarrow Aux NP VP$

- WH Questions
When did the plane leave?
 $S \rightarrow NP_{WH} Aux NP VP$

Grammar Fragments: NP

- Each NP has a central critical noun called **head**
- The head of an NP can be expressed using
 - Pre-nominals: the words that can come before the head
 - Post-nominals: the words that can come after the head

Grammar Fragments: NP

■ Pre-nominals

- Simple lexical items: *the, this, a, an, ...*
a car
- Simple possessives
John's car
- Complex recursive possessives
John's sister's friend's car
- Quantifiers, cardinals, ordinals...
three cars
- Adjectives
large cars

Grammar Fragments: NP

18

- Post-nominals
 - Prepositional phrases
flight from Seattle
 - Non-finite clauses
flight arriving before noon
 - Relative clauses
flight that serves breakfast

Agreement

- Having constraints that hold among various constituents
- Considering these constraints in a rule or set of rules

Example: determiners and the head nouns in NPs have to agree in number

This flight ✓

Those flights ✓

This flights ✗

Those flight ✗

- Grammars that do not consider constraints will **over-generate**
 - Accepting and assigning correct structures to grammatical examples (*this flight*)
 - But also accepting incorrect examples (*these flight*)

Agreement at sentence level

- Considering similar constraints at sentence level

Example: subject and verb in sentences have to agree in number and person

John flies ✓

We fly ✓

John fly ✗

We flies ✗

Agreement

21

- Possible CFG solution

$$S_{sg} \rightarrow NP_{sg} VP_{sg}$$
$$S_{pl} \rightarrow NP_{pl} VP_{pl}$$
$$NP_{sg} \rightarrow Det_{sg} N_{sg}$$
$$NP_{pl} \rightarrow Det_{pl} N_{pl}$$
$$VP_{sg} \rightarrow V_{sg} NP_{sg}$$
$$VP_{pl} \rightarrow V_{pl} NP_{pl}$$

...

- Shortcoming:

- Introducing many rules in the system

Grammar Fragments: VP

22

- VPs consist of a head verb along with zero or more constituents called **arguments**

$VP \rightarrow V$	<i>disappear</i>
$VP \rightarrow V NP$	<i>prefer a morning flight</i>
$VP \rightarrow V PP$	<i>fly on Thursday</i>
$VP \rightarrow V NP PP$	<i>leave Boston in the morning</i>
$VP \rightarrow V NP NP$	<i>give me the flight number</i>

- Arguments
 - Obligatory: complement
 - Optional: adjunct

Sub-categorization

23

- Even though there are many valid VP rules, not all verbs are allowed to participate in all VP rules

disappear a morning flight ✗

- Solution:
 - Subcategorizing the verbs according to the sets of VP rules that they can participate in
 - This is a modern take on the traditional notion of transitive/intransitive
 - Modern grammars may have 100s or such classes

Sub-categorization

24

- Example:

Sneeze	<i>John sneezed</i>
Find	<i>Please find [a flight to NY]_{NP}</i>
Give	<i>Give [me]_{NP}[a cheaper fair]_{NP}</i>
Help	<i>Can you help [me]_{NP}[with a flight]_{PP}</i>
Prefer	<i>I prefer [to leave earlier]_{TO-VP}</i>
Told	<i>I was told [United has a flight]_s</i>

John sneezed the book ✗
I prefer United has a flight ✗
Give with a flight ✗

Sub-categorization

25

- The over-generation problem also exists in VP rules
 - Permitting the presence of strings containing verbs and arguments that do not go together

John sneezed the book

$VP \rightarrow V NP$

- Solution:
 - Similar to agreement phenomena, we need a way to formally express the constraints

Parsing Algorithms

26

■ Top-Down

- Starting with the rules that give us an S, since trees should be rooted with an S
- Working on the way down from S to the words

■ Bottom-Up

- Starting with trees that link up with the words, since trees should cover the input words
- Working on the way up from words to larger and larger trees

Top-Down vs. Bottom-Up

27

■ 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

Top-Down vs. Bottom-Up

28

- In both cases, we left out how to keep track of the search space and how to make choices

- Solutions
 - Backtracking
 - Making a choice, if it works out then fine
 - If not, then back up and make a different choice
⇒ duplicated work

 - Dynamic programming
 - Avoiding repeated work
 - Solving exponential problems in polynomial time
 - Storing ambiguous structures efficiently

- CKY: bottom-up
- Early: top-down

Outline

30

- ① Phrase Structure
- ② Syntactic Parsing
CKY Algorithm
- ③ Statistical Parsing

Chomsky Normal Form

31

- Each grammar can be represented by a set of binary rules

$$A \rightarrow B C$$

$$A \rightarrow w$$

A, B, C are non-terminals w is a terminal

Chomsky Normal Form

32

- Converting to Chomsky normal form

$$A \rightarrow B C D$$

$$X \rightarrow B C$$

$$A \rightarrow X D$$

X does not occur anywhere else in the the grammar

Chomsky Normal Form

33

- Converting to Chomsky normal form

$$A \rightarrow B$$

$$B \rightarrow C D$$

$$A \rightarrow C D$$

CKY Parsing

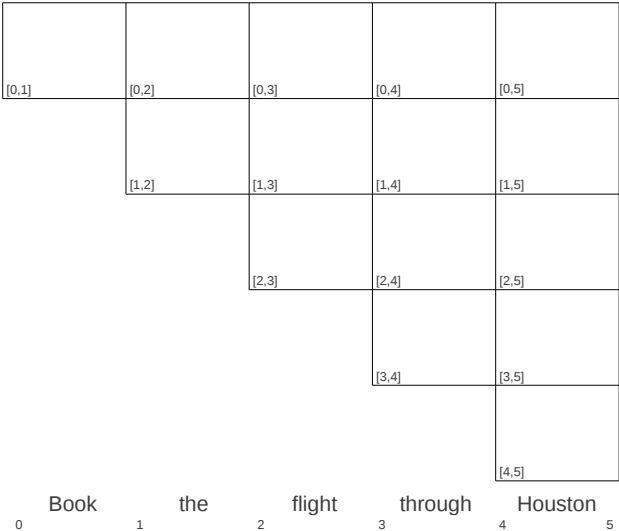
34

$$A \rightarrow B C$$

- 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 a k such that $i < k < j$
 - B spans from i to k
 - C spans from k to j

CKY Parsing

35



CKY Parsing

36

N → book _[0,1] V → book _[0,1]						
[0,1]	[0,2]	[0,3]	[0,4]	[0,5]		
	Det → the _[1,2]					
	[1,2]	[1,3]	[1,4]	[1,5]		
		N → flight _[2,3]				
		[2,3]	[2,4]	[2,5]		
			Prep → through _[3,4]			
			[3,4]	[3,5]		
				N → houston _[4,5]		
				[4,5]		
	Book	the	flight	through	Houston	
	0	1	2	3	4	5

CKY Parsing

37

N → book _[0,1] V → book _[0,1] NP → N _[0,1] VP → V _[0,1] S → VP _[0,1]				
[0,1]	[0,2]	[0,3]	[0,4]	[0,5]
	Det → the _[1,2]			
	[1,2]	[1,3]	[1,4]	[1,5]
		N → flight _[2,3] NP → N _[2,3]		
		[2,3]	[2,4]	[2,5]
			Prep → through _[3,4]	
			[3,4]	[3,5]
				N → houston _[4,5] NP → N _[4,5]
				[4,5]

0 Book 1 the 2 flight 3 through 4 Houston 5

CKY Parsing

38

N → book _[0,1] V → book _[0,1] NP → N _[0,1] VP → V _[0,1] S → VP _[0,1] [0,1]				
	Det → the _[1,2] [1,2]	NP → Det _[1,2] , N _[2,3] [1,3]		
		N → flight _[2,3] NP → N _[2,3] [2,3]		
			Prep → through _[3,4] [3,4]	
				N → houston _[4,5] NP → N _[4,5] [4,5]

0 Book 1 the 2 flight 3 through 4 Houston 5

CKY Parsing

39

N → book _[0,1] V → book _[0,1] NP → N _[0,1] VP → V _[0,1] S → VP _[0,1] [0,1]				
	Det → the _[1,2] [1,2]	NP → Det _[1,2] , N _[2,3] [1,3]		
		N → flight _[2,3] NP → N _[2,3] [2,3]		
			Prep → through _[3,4] [3,4]	PP → Prep _[3,4] , NP _[4,5] [3,5]
				N → houston _[4,5] NP → N _[4,5] [4,5]

0 Book 1 the 2 flight 3 through 4 Houston 5

CKY Parsing

40

N → book _[0,1] V → book _[0,1] NP → N _[0,1] VP → V _[0,1] S → VP _[0,1] [0,1]		NP → NP _[0,1] , NP _[1,3] VP → VP _[0,1] , NP _[1,3] S → VP _[0,3] [0,3]		
	Det → the _[1,2] [1,2]	NP → Det _[1,2] , N _[2,3] [1,3]		
		N → flight _[2,3] NP → N _[2,3] [2,3]		
			Prep → through _[3,4] [3,4]	PP → Prep _[3,4] , NP _[4,5] [3,5]
				N → houston _[4,5] NP → N _[4,5] [4,5]

0 Book 1 the 2 flight 3 through 4 Houston 5

CKY Parsing

41

N → book _[0,1] V → book _[0,1] NP → N _[0,1] VP → V _[0,1] S → VP _[0,1] [0,1]		NP → NP _[0,1] , NP _[1,3] VP → VP _[0,1] , NP _[1,3] S → VP _[0,3] [0,3]		
	Det → the _[1,2] [1,2]	NP → Det _[1,2] , N _[2,3] [1,3]		
		N → flight _[2,3] NP → N _[2,3] [2,3]		NP → NP _[2,3] , PP _[3,5] [2,5]
			Prep → through _[3,4] [3,4]	PP → Prep _[3,4] , NP _[4,5] [3,5]
				N → houston _[4,5] NP → N _[4,5] [4,5]

0 Book 1 the 2 flight 3 through 4 Houston 5

CKY Parsing

42

N → book _[0,1] V → book _[0,1] NP → N _[0,1] VP → V _[0,1] S → VP _[0,1] [0,1]		NP → NP _[0,1] , NP _[1,3] VP → VP _[0,1] , NP _[1,3] S → VP _[0,3] [0,3]		
	Det → the _[1,2] [1,2]	NP → Det _[1,2] , N _[2,3] [1,3]		NP → NP _[1,3] , PP _[3,5] [1,5]
		N → flight _[2,3] NP → N _[2,3] [2,3]		NP → NP _[2,3] , PP _[3,5] [2,5]
			Prep → through _[3,4] [3,4]	PP → Prep _[3,4] , NP _[4,5] [3,5]
				N → houston _[4,5] NP → N _[4,5] [4,5]

0 Book 1 the 2 flight 3 through 4 Houston 5

CKY Parsing

43

N → book _[0,1] V → book _[0,1] NP → N _[0,1] VP → V _[0,1] S → VP _[0,1] [0,1]		NP → NP _[0,1] , NP _[1,3] VP → VP _[0,1] , NP _[1,3] S → VP _[0,3] [0,3]		VP → VP _[0,1] , NP _[1,5] VP' → VP _[0,3] , PP _[3,5] S → VP _[0,5] S → VP' _[0,5] [0,5]
	Det → the _[1,2] [1,2]	NP → Det _[1,2] , N _[2,3] [1,3]		NP → NP _[1,3] , PP _[3,5] [1,5]
		N → flight _[2,3] NP → N _[2,3] [2,3]		NP → NP _[2,3] , PP _[3,5] [2,5]
			Prep → through _[3,4] [3,4]	PP → Prep _[3,4] , NP _[4,5] [3,5]
				N → houston _[4,5] NP → N _[4,5] [4,5]

ambiguity

0 Book 1 the 2 flight 3 through 4 Houston 5

Outline

44

- 1 Phrase Structure
- 2 Syntactic Parsing
CKY Algorithm
- 3 Statistical Parsing

- Grammar G consists of
 - Terminals (T)
 - Non-terminals (N)
 - Start symbol (S)
 - Rules (R)
 - Probability function (P)
 - $P: R \rightarrow [0, 1]$
 - $\forall X \in N, \sum_{X \rightarrow \lambda \in R} P(X \rightarrow \lambda) = 1$

CFG

46

 $S \rightarrow NP VP$ $S \rightarrow VP$ $NP \rightarrow N$ $NP \rightarrow Det N$ $NP \rightarrow NP NP$ $NP \rightarrow NP PP$ $VP \rightarrow V$ $VP \rightarrow VP PP$ $VP \rightarrow VP NP$ $PP \rightarrow Prep NP$ $N \rightarrow \text{book}$ $V \rightarrow \text{book}$ $Det \rightarrow \text{the}$ $N \rightarrow \text{flight}$ $Prep \rightarrow \text{through}$ $N \rightarrow \text{Houston}$

PCFG

47

$S \rightarrow NP VP$	0.9		
$S \rightarrow VP$	0.1		
$NP \rightarrow N$	0.3	$N \rightarrow \text{book}$	0.5
$NP \rightarrow Det N$	0.4	$V \rightarrow \text{book}$	1.0
$NP \rightarrow NP NP$	0.1	$Det \rightarrow \text{the}$	1.0
$NP \rightarrow NP PP$	0.2	$N \rightarrow \text{flight}$	0.4
$VP \rightarrow V$	0.1	$Prep \rightarrow \text{through}$	1.0
$VP \rightarrow VP PP$	0.3	$N \rightarrow \text{Houston}$	0.1
$VP \rightarrow VP NP$	0.6		
$PP \rightarrow Prep NP$	1.0		

Treebank

48

- A treebank is a corpus in which each sentence has been paired with a parse tree
- These are generally created by
 - Parsing the collection with an automatic parser
 - Correcting each parse by human annotators if required
- Requirement:
detailed annotation guidelines that provide
 - A POS tagset
 - A grammar
 - Annotation schema
 - Instructions for how to deal with particular grammatical constructions

Penn Treebank

49

- Penn Treebank is a widely used treebank for English
 - Most well-known section: Wall Street Journal Section
 - 1 M words from 1987-1989

```
(S (NP (NNP John))
  (VP (VPZ flies)
    (PP (IN to)
      (NNP Paris))))
(. .))
```

Statistical Parsing

50

- Considering the corresponding probabilities while parsing a sentence
- Selecting the parse tree which has the highest probability
- Tree and string probabilities
 - $P(t)$: the probability of a tree t
 - Product of the probabilities of the rules used to generate the tree
 - $P(s)$: the probability of a string s
 - Sum of the probabilities of the trees which created to parse the string

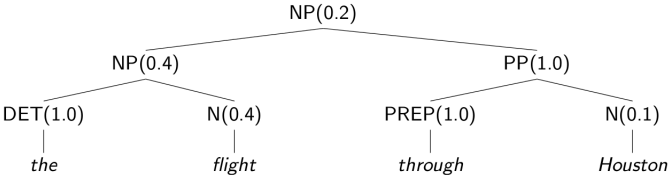
PCFG

51

$S \rightarrow NP VP$	0.9		
$S \rightarrow VP$	0.1		
$NP \rightarrow N$	0.3	$N \rightarrow \text{book}$	0.5
$NP \rightarrow Det N$	0.4	$V \rightarrow \text{book}$	1.0
$NP \rightarrow NP NP$	0.1	$Det \rightarrow \text{the}$	1.0
$NP \rightarrow NP PP$	0.2	$N \rightarrow \text{flight}$	0.4
$VP \rightarrow V$	0.1	$Prep \rightarrow \text{through}$	1.0
$VP \rightarrow VP PP$	0.3	$N \rightarrow \text{Houston}$	0.1
$VP \rightarrow VP NP$	0.6		
$PP \rightarrow Prep NP$	1.0		

Statistical Parsing

52



$$P(t) = 0.2 \times 0.4 \times 1.0 \times 1.0 \times 0.4 \times 1.0 \times 0.1 = 0.0032$$

Probabilistic CKY Parsing

53

0.5
1.0
0.3*0.5=0.15
0.1*1.0=0.1
0.1*0.1=0.01

<p>N → book_[0,1] V → book_[0,1] NP → N_[0,1] VP → V_[0,1] S → VP_[0,1]</p> <p>[0,1]</p>	<p>$0.1 \cdot 0.15 \cdot 0.16 = 0.0024$ $0.6 \cdot 0.1 \cdot 0.16 = 0.0096$ $0.1 \cdot 0.0096 = 0.00096$</p>	<p>NP → NP_[0,1], NP_[1,3] VP → VP_[0,1], NP_[1,3] S → VP_[0,3]</p> <p>[0,3]</p>	<p>[0,4]</p>	<p>VP → VP_[0,1], NP_[1,5] VP' → VP_[0,3], PP_[3,5] S → VP_[0,5] S → VP'_[0,5]</p> <p>[0,5]</p>
<p>1.0</p>	<p>Det → the_[1,2]</p> <p>[1,2]</p>	<p>NP → Det_[1,2], N_[2,3] $0.4 \cdot 1.0 \cdot 0.4 = 0.16$</p> <p>[1,3]</p>	<p>[1,4]</p>	<p>NP → NP_[1,3], PP_[3,5]</p> <p>[1,5]</p>
	<p>0.4 $0.3 \cdot 0.4 = 0.12$</p>	<p>N → flight_[2,3] NP → N_[2,3]</p> <p>[2,3]</p>	<p>[2,4]</p>	<p>NP → NP_[2,3], PP_[3,5]</p> <p>[2,5]</p>
			<p>Prep → through_[3,4]</p> <p>[3,4]</p>	<p>PP → Prep_[3,4], NP_[4,5]</p> <p>[3,5]</p>
				<p>N → houston_[4,5] NP → N_[4,5]</p> <p>[4,5]</p>

0 Book 1 the 2 flight 3 through 4 Houston 5

Exercise

54

- Implement the probabilistic CKY algorithm which works based on the grammar rules R .

Further Reading

- Speech and Language Processing
 - Chapters 12, 13, 14, 15