

# Statistical Parsing Chapter 14

Lecture #9

October 2012

1

## Parse Disambiguation

- In the previous chapter we have seen several instances of parsing ambiguity: coordination ambiguity and attachment ambiguity
- So far – we return every parse and let later modules deal with the ambiguity
- Can we use probabilistic methods to choose most likely interpretation?

2

## Probabilistic Parsing

- Probabilistic Context Free Grammar: a probabilistic grammar which favors more common rules
- Augment each rule with its associated probability
- Modify parser so that it returns most likely parse (CKY Algorithm)
- Problems and augmentations to the basic model

3

## Probability Model

- Attach probabilities to grammar rules – representing the probability that a given non-terminal on the rule's LHS will be expanded to the sequence on the rule's LHS.
- The expansions for a given non-terminal sum to 1
  - VP → verb .55
  - VP → verb NP .40
  - VP → verb NP NP .05
- Probability captures  $P(\text{RHS} | \text{LHS})$

4

Figure 14.1

Grammar		Lexicon
$S \rightarrow NP VP$	.80	$Det \rightarrow that [.10]   a [.30]   the [.60]$
$S \rightarrow Aux NP VP$	.15	$Noun \rightarrow book [.10]   flight [.30]$
$S \rightarrow VP$	.05	$  meal [.15]   money [.05]$
$NP \rightarrow Pronoun$	.35	$  flights [.40]   dinner [.10]$
$NP \rightarrow Proper-Noun$	.30	$Verb \rightarrow book [.30]   include [.30]$
$NP \rightarrow Det Nominal$	.20	$  prefer; [.40]$
$NP \rightarrow Nominal$	.15	$Pronoun \rightarrow I [.40]   she [.05]$
$Nominal \rightarrow Noun$	.75	$  me [.15]   you [.40]$
$Nominal \rightarrow Nominal Noun$	.20	$Proper-Noun \rightarrow Houston [.60]$
$Nominal \rightarrow Nominal PP$	.05	$  NWA [.40]$
$VP \rightarrow Verb$	.35	$Aux \rightarrow does [.60]   can [.40]$
$VP \rightarrow Verb NP$	.20	$Preposition \rightarrow from [.30]   to [.30]$
$VP \rightarrow Verb NP PP$	.10	$  on [.20]   near [.15]$
$VP \rightarrow Verb PP$	.15	$  through [.05]$
$VP \rightarrow Verb NP NP$	.05	
$VP \rightarrow VP PP$	.15	
$PP \rightarrow Preposition NP$	1.0	

## Using a PCFS

- A PCFS can be used to estimate a number of useful probabilities concerning a sentence and its parse trees
  - Probability of a particular parse tree (useful for disambiguation)
  - Probability of a sentence or piece of a sentence (useful for language modeling)

How?

- A derivation (tree) consists of the set of grammar rules that are in the tree
- The probability of a tree is just the product of the probabilities of the rules in the derivation.

6

Figure 14.2

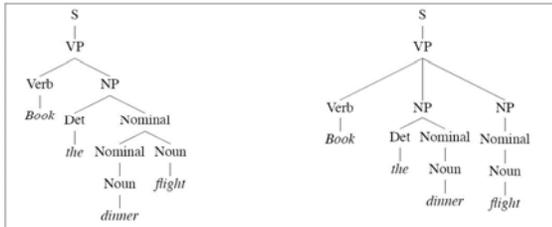


Figure 14.2 continued

Rules	P	Rules	P
S → VP	.05	S → VP	.05
VP → Verb NP	.20	VP → Verb NP NP	.10
NP → Det Nominal	.20	NP → Det Nominal	.20
Nominal → Nominal Noun	.20	NP → Nominal	.15
Nominal → Noun	.75	Nominal → Noun	.75
		Nominal → Noun	.75
Verb → book	.30	Verb → book	.30
Det → the	.60	Det → the	.60
Noun → dinner	.10	Noun → dinner	.10
Noun → flights	.40	Noun → flights	.40

## Probability Model

$P(T, S) = P(T)P(S | T) = P(T)$ ; since  $P(S | T) = 1$

$$P(T, S) = P(T) = \prod_{n \in T} p(r_n)$$

$$P(T_{left}) = .05 * .20 * .20 * .20 * .75 * .30 * .60 * .10 * .40 = 2.2 \times 10^{-6}$$

$$P(T_{right}) = .05 * .10 * .20 * .15 * .75 * .75 * .30 * .60 * .10 * .40 = 6.1 \times 10^{-7}$$

9

## Probability Model

- The probability of a word sequence P(S) is the probability of its tree in the unambiguous case (i.e., where there is exactly one tree).
- In the case where there is ambiguity (multiple trees) the probability of the sequence is the sum of the probabilities of the trees.

10

## Parsing to get most likely parse

- Can do with a simple extension of our parsing algorithms – book does CKY (and indicates that is most used version).
- Essentially – give each constituent that is in the table a probability (they refer to this as another dimension in the table), when a new constituent, C, is found to be added to the table at cell [I, J], only add it if that cell either does not contain a constituent C or if the probability of this new constituent is less than the probability of the existing one (in which case, you overwrite the old one).
- Assuming rule is  $C \rightarrow c_1 c_2$
- New prob = prob(rule) x prob(c1) x prob(c2)

11

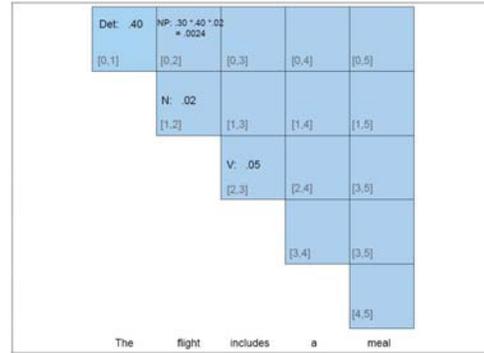
Figure 14.3

```

function PROBABILISTIC-CKY(words, grammar) returns most probable parse
and its probability
for j ← from 1 to LENGTH(words) do
  for all { A | A → words[j] ∈ grammar }
    table[j-1, j, A] ← P(A → words[j])
  for i ← from j-2 down to 0 do
    for all { A | A → BC ∈ grammar,
              and table[i, k, B] > 0 and table[k, j, C] > 0 }
      if (table[i, j, A] < P(A → BC) × table[i, k, B] × table[k, j, C]) then
        table[i, j, A] ← P(A → BC) × table[i, k, B] × table[k, j, C]
        back[i, j, A] ← {k, B, C}
  return BUILD-TREE(back[1, LENGTH(words), S], table[1, LENGTH(words), S])
    
```

$S \rightarrow NP VP$	.80	$Det \rightarrow the$	.40
$NP \rightarrow Det N$	.30	$Det \rightarrow a$	.40
$VP \rightarrow V NP$	.20	$N \rightarrow meal$	.01
$V \rightarrow includes$	.05	$N \rightarrow flight$	.02

Figure 14.4



## Learning PCFG Rule Probabilities

- Learn from a treebank, a corpus of already parsed sentences
  - So for example, to get the probability of a particular VP rule just count all the times the rule is used and divide by the number of VPs overall
- Count by parsing with a non-probabilistic parser.
  - Issue is ambiguity – inside-out-algorithm - sort of boot-strap – parse a sentence, compute a probability for each parse, use these probabilities to weight the counts, re-estimate the rule probabilities, and so on, until our probabilities converge.

15

## Problems with PCFGs

- Probability model we are using is just based on the rules in the derivation... and these are context free rules
  - Poor independence assumptions miss structural dependencies between rules since cannot take into account in the derivation a rule is used.
  - Lack of sensitivity to lexical dependencies
    - Do have probability associated with N-> bank
    - But verb subcategorization and prepositional phrase attachment might depend on the particular words being used.
    - Use lexical heads as part of rule – then you run into problems with sparse data – so need to make independence assumptions to reduce amount of data needed.

16

## Structural Dependencies between Rules

$NP \rightarrow det NN$  .28  
 $NP \rightarrow Pronoun$  .25

- These probabilities should depend on **where** the NP is being used:

	Pronoun	Non-Pronoun
Subject	91%	9%
Object	34%	66%

- SolutionSplit non-terminal into 2 (e.g., using parent annotation  $NP^A S$  vs  $NP^B VP$ ) and learn rule probabilities for split rules.

17

Figure 14.8

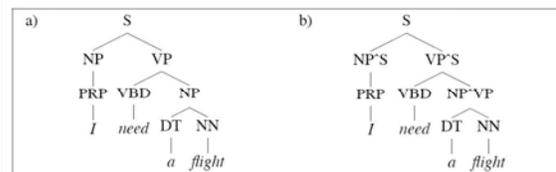


Figure 14.5

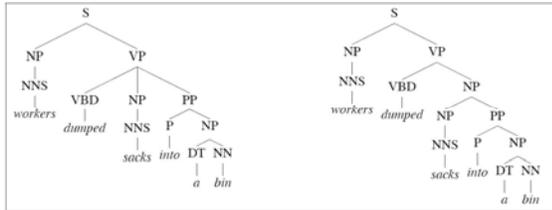
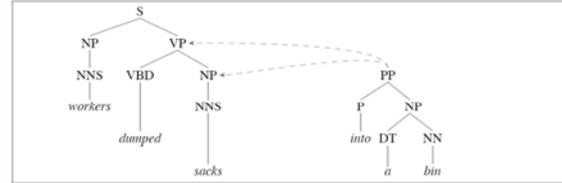


Figure 14.6 – KFM added to slide

*Note: if we prefer the rule for VP attachment then it will be incorrect for sentences like “fishermen caught tons of herring”*

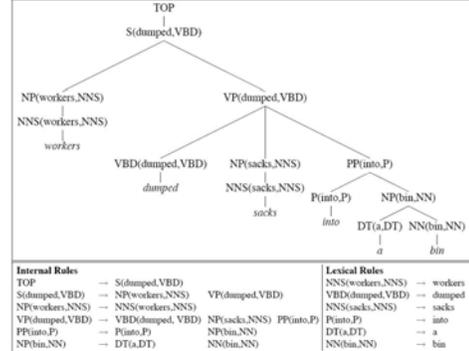


*Seems lexical in nature: affinity between the verb “dumped” and the preposition “into” is greater than the affinity between the noun “sacks” and the preposition “into”*

## Lexical Dependencies

- Must add lexical dependencies to the scheme and condition the rule probabilities on the actual words
- What words?
  - Make use of the notion of the **head** of a phrase
  - The head is intuitively the most important lexical item in the phrase – and there are some rules for identifying
    - Head of an NP is its noun
    - Head of a VP is its verb
    - Head of a sentence comes from its VP
    - Head of a PP is its preposition
  - Use a lexicalized grammar in which each non-terminal in the tree is annotated with its lexical head

Figure 14.10



## Issues with Learning

- Not likely to have significant counts in any treebank to actually learn these probabilities.
- Solution: Make as many independence assumptions as you can and learn from these
- Different modern parsers make different independence assumptions – E.g., Collins parser have head and dependents on left are assumed independent of each other and independent of dependents on right (which make similar assumptions about the left dependents)

## Summary

- Probabilistic Context-Free Grammars
  - Help us deal with ambiguity by preferring more likely parses
- Grammar rules have attached probabilities which capture the probability of the rule's RHS given its LHS (probabilities of all rules with same LHS sum to 1)
- We can compute the probability of a tree (product of the probabilities of the rules used)
- Can parse using augmented algorithms
- Can learn probabilities from a tree bank
- PCFGs have problems with independence assumptions and with lack of lexical conditioning
- Some solutions exist – problems of data sparsity