

FNLP Tutorial 5

1 Attachment Ambiguity in Trees

A common source of attachment ambiguity in English comes from prepositional phrases. The relevant grammar rules include:

$VP \rightarrow V \ NP$

$VP \rightarrow VP \ PP$

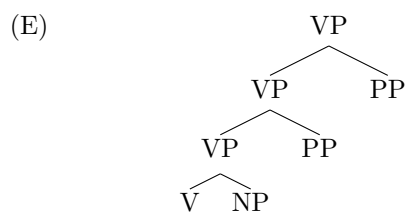
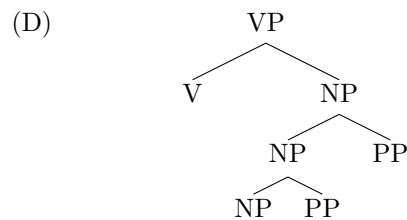
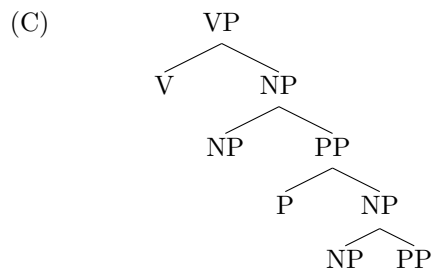
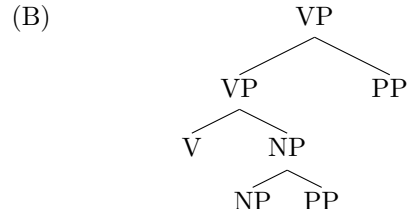
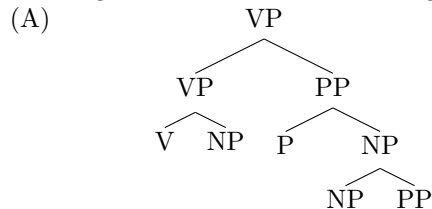
$NP \rightarrow NP \ PP$

$PP \rightarrow P \ NP$

Here are five verb phrases:

1. read the paper about the Turing test on my laptop
2. read the paper about the success of Transformer
3. read the paper by Alan Turing about the Turing test
4. joined the meeting in the office with the blue door
5. joined the online meeting from my office via Zoom

Pictured below are five partial trees. Match the phrases to the trees which best capture their meanings. Is there a source of ambiguity with the first sentence?



Solution 1B; 2C; 3D; 4A; 5E The first sentence could also be C, if we understood that the paper was about a Turing test done on my laptop. However, this is semantically implausible -we know Turing tests are not done on laptops-.

2 Writing a context-free grammar

Consider this small corpus of English noun phrases.

1. this book
2. the great green dragon
3. these green leaves
4. the leaves in the hallway

- Write a context-free grammar that generates at least the above noun phrases. Use the following part-of-speech tags as pre-terminal symbols: determiner (DT), adjective (A), preposition (P) and noun (N).
- Show a derivation tree of a grammatical noun phrase that is not included in the corpus. Does your grammar overgenerate? That is, does it generate sequences that are not grammatical English noun phrases? If so, refine your grammar to fix the issue that you mentioned.

Solution We will start with the following set of POS tags:

DT determiner, e.g. this and these

A adjective, e.g. great and green

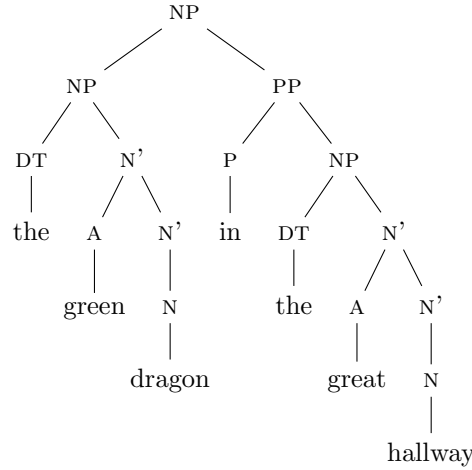
P preposition, e.g. in

N noun, e.g. book

Our first attempt might look something like this (NP is the start symbol). We write $X \rightarrow Y|Z$ as a shorthand for the two productions $X \rightarrow Y$ and $X \rightarrow Z$.

$$\begin{aligned} \text{NP} &\rightarrow \text{DT} \text{ N}' \\ \text{NP} &\rightarrow \text{NP} \text{ PP} \\ \text{N}' &\rightarrow \text{A} \text{ N}' \\ \text{N}' &\rightarrow \text{N} \\ \text{PP} &\rightarrow \text{P} \text{ NP} \\ \text{DT} &\rightarrow \text{the}|\text{this}|\text{these} \\ \text{P} &\rightarrow \text{in} \\ \text{A} &\rightarrow \text{great}|\text{green} \\ \text{N} &\rightarrow \text{book}|\text{dragon}|\text{leaves}|\text{hallway} \end{aligned}$$

We can derive the noun phrase "the green dragon in the great hallway" from this grammar:



Unfortunately, this grammar also generates ungrammatical noun phrases like:

- *this leaves
- *these great book

To prevent this from happening, our grammar must encode number information (if something is singular or plural). In particular, we need to encode this information for determiners and nouns and make sure that only singular determiners and singular nouns can combine (and analogously for plural).

We first refine our set of POS tags slightly to encode number information, i.e. we introduce DT_{SG} , DT_{PL} and N_{SG} and N_{PL} . We also modify production rules to track this information:

$$\begin{aligned}
 NP &\rightarrow DT_{SG} \quad N'_{SG} \\
 NP &\rightarrow DT_{PL} \quad N'_{PL} \\
 NP &\rightarrow NP \quad PP \\
 N'_{SG} &\rightarrow A \quad N'_{SG} \\
 N'_{PL} &\rightarrow A \quad N'_{PL} \\
 N'_{SG} &\rightarrow N_{SG} \\
 N'_{PL} &\rightarrow N_{PL} \\
 PP &\rightarrow P \quad NP \\
 DT_{SG} &\rightarrow \text{the}|\text{this} \\
 DT_{PL} &\rightarrow \text{these} \\
 P &\rightarrow \text{in} \\
 A &\rightarrow \text{great}|\text{green} \\
 N_{SG} &\rightarrow \text{book}|\text{dragon}|\text{hallway} \\
 N_{PL} &\rightarrow \text{leaves}
 \end{aligned}$$

This grammar still generates noun phrases that are disfavoured if not ungrammatical, like:

- ??the green great dragon

3 The CKY Algorithm

1. Determine the potential parses of the sentence “The women fish with bait” according to the PCFG below, using the CKY algorithm. Number the symbols you put in the chart in the

order they would be computed. For each chart item, assume that candidate grammar rules are searched in the order they are listed below. Use these numbers to construct backpointers in the chart. For example, when constructing the 6th symbol NP out of a determiner with symbol number 2 and a noun with symbol number 3, you could indicate that by putting '6: NP (2, 3)'. (This is not a widely used custom, but a useful one when doing CKY on paper for the context of this tutorial.)

2. What are the probabilities of the parse trees?

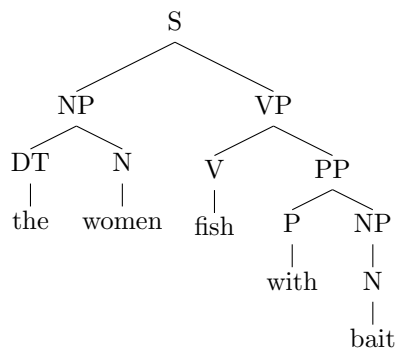
$S \rightarrow NP \quad VP \quad (1)$
 $NP \rightarrow DT \quad N \quad (0.5)$
 $NP \rightarrow N \quad N \quad (0.1)$
 $NP \rightarrow N \quad (0.4)$
 $VP \rightarrow V \quad PP \quad (1)$
 $PP \rightarrow P \quad NP \quad (1)$
 $DT \rightarrow \text{the} \quad (1)$
 $P \rightarrow \text{with} \quad (1)$
 $V \rightarrow \text{fish} \quad (1)$
 $N \rightarrow \text{women} \quad (0.5) \mid \text{fish} \quad (0.3) \mid \text{bait} \quad (0.2)$

Solution

1. Here is the CKY chart:

the	women	fish	with	bait
1: DT p=1	10: NP(1,2) p=0.25 2: N p=0.5, 3: NP(2) p=0.2	11: NP(2,5) p=0.015 4: V p=1, 5: N p=0.3, 6: NP(5) p=0.12	7: P p=1	15: S(10,13) p=0.02 14: S(3,13) p=0.016 13: VP(4,12) p=0.08 12: PP(7,9) p=0.08 8: N p=0.2, 9: NP(8) p=0.08

The resulting parse looks like this:



2. There is only one valid parse. The parse tree has the following probability: $1 \times 0.5 \times 1 \times 1 \times 0.2 \times 0.5 \times 0.4 \times 1 \times 1 \times 1 = 0.02$.