

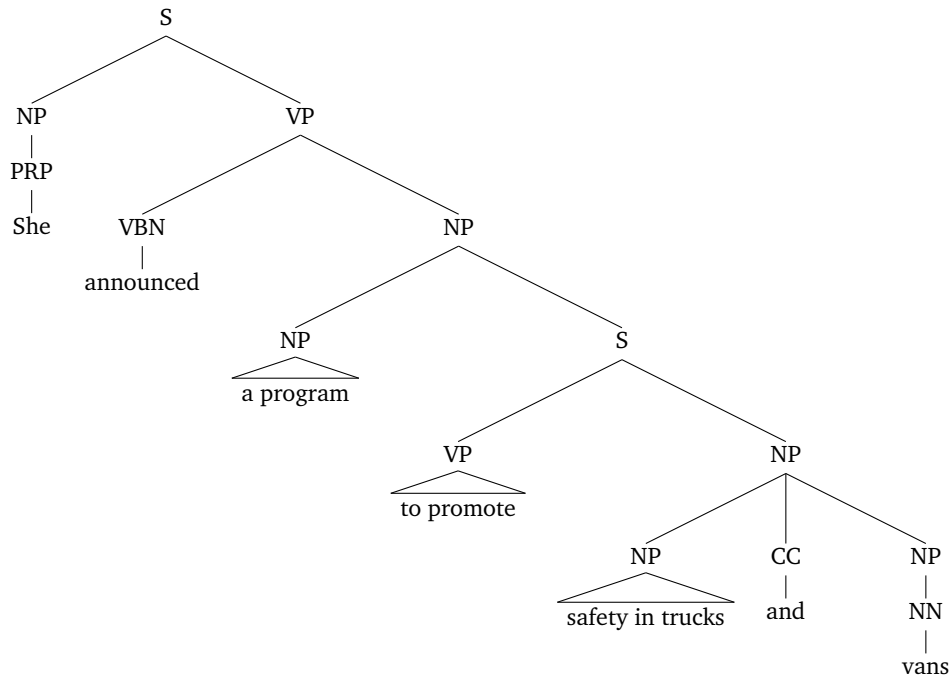
Mining Unstructured Data

Exercises on Constituent Parsing

Context Free Grammars

Exercise 1.

Consider the sentence *She announced a program to promote safety in trucks and vans* and the following syntactic tree of one of its possible interpretations, in which the program promotes safety in trucks, and also promotes vans:

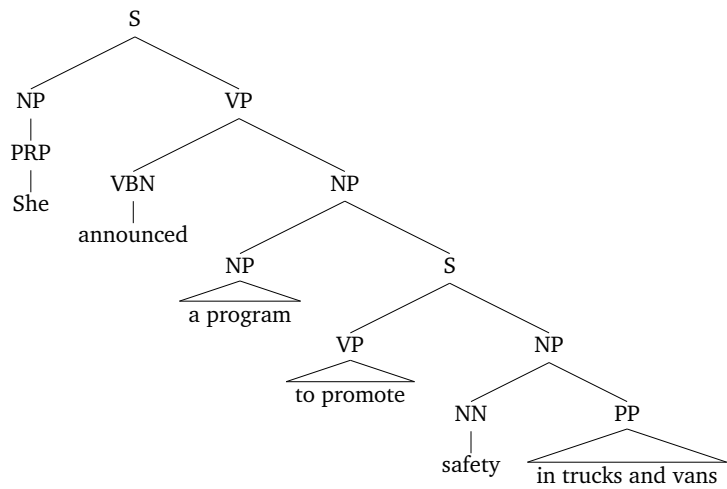


1. Draw the trees for at least three other interpretations for this sentence
2. Draw the trees for at least two interpretations for each of the following sentences
 - *The post office will hold out discounts and service concessions as incentives*
 - *They are hunting lions and tigers*
 - *Monty flies like mosquitoes*

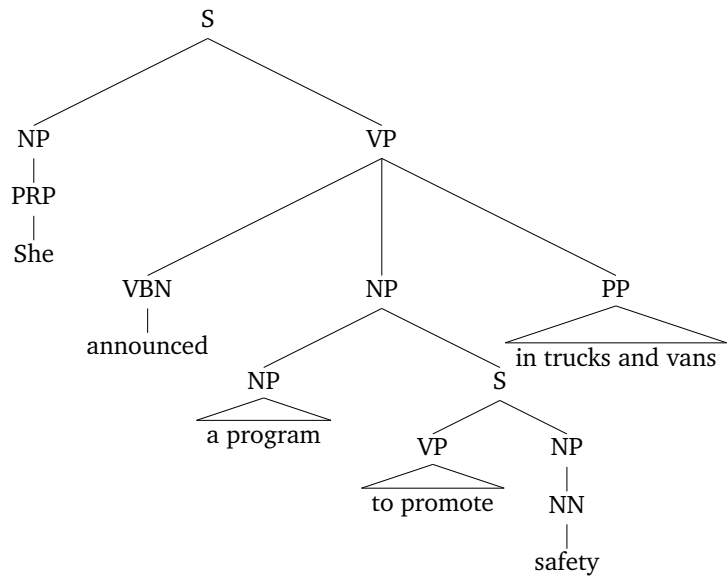
SOLUTION

- Find three interpretations:

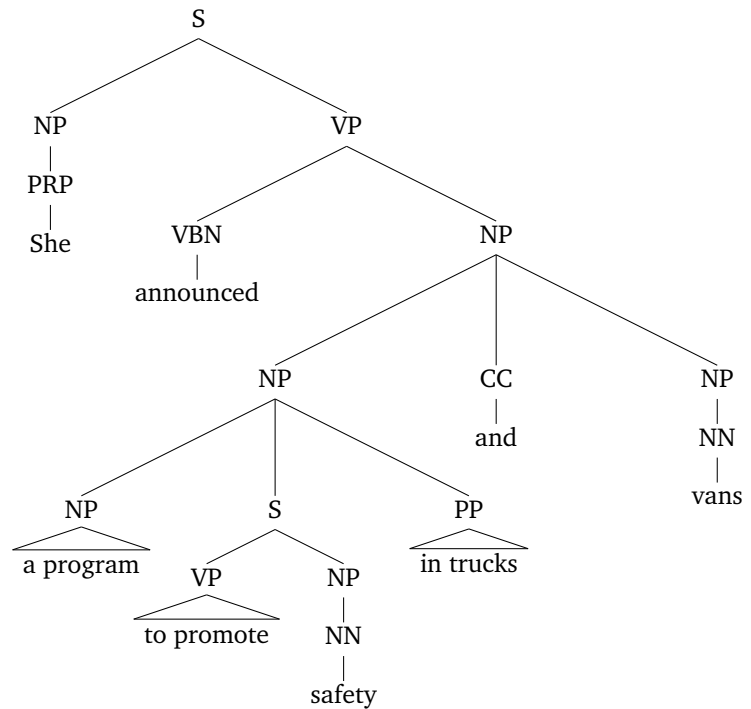
Interpretation 1: The announced program promotes safety in both trucks and vans.



Interpretation 2: The program is announced in trucks and vans.

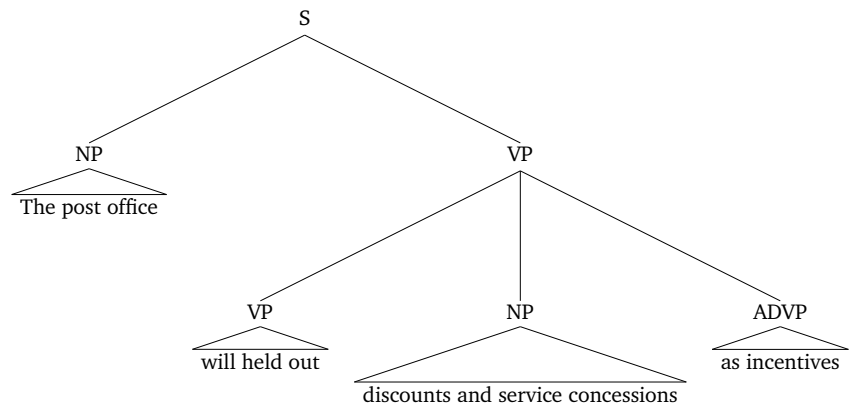


Interpretation 3: The announced program promotes safety in trucks. Vans are also announced.

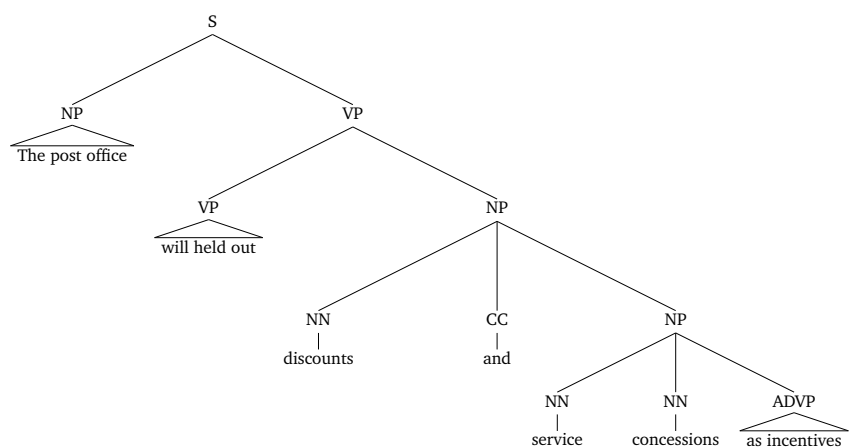


2. Find two interpretations for each sentence

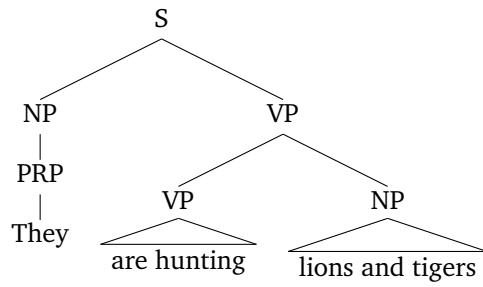
Sentence 1, Interpretation 1:
Discounts and concessions are held out as incentives



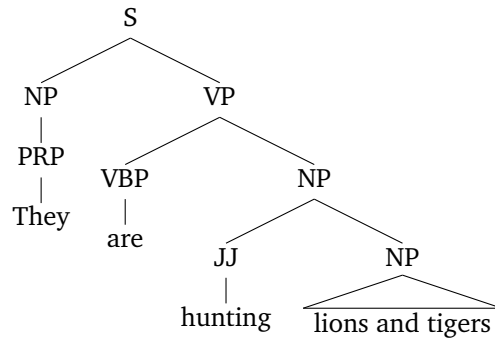
Sentence 1, Interpretation 2:
Discounts and concessions are held out. Concessions look like incentives



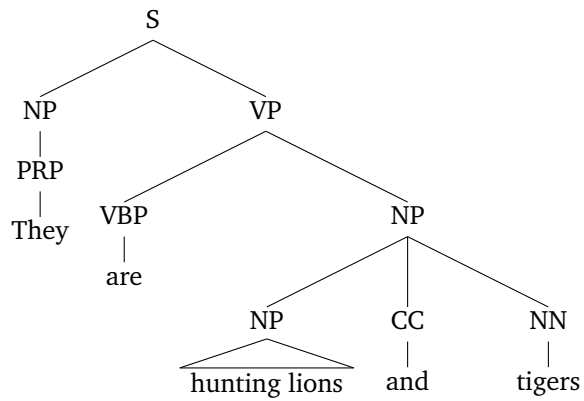
Sentence 2, Interpretation 1:
Someone is hunting big felines.



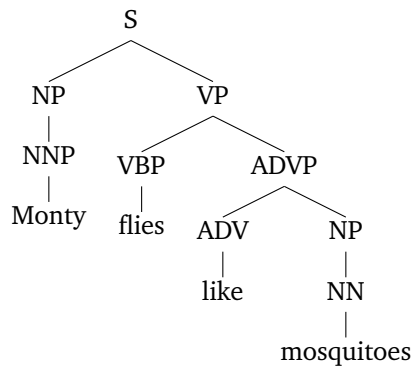
Sentence 2, Interpretation 2:
Those animals are big felines that hunt.



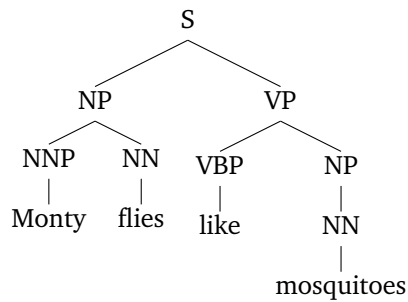
Sentence 2, Interpretation 3:
Those animals are lions that hunt, and also tigers



Sentence 3, Interpretation 1:
Someone named Monty moves through the air in the same way than mosquitoes do

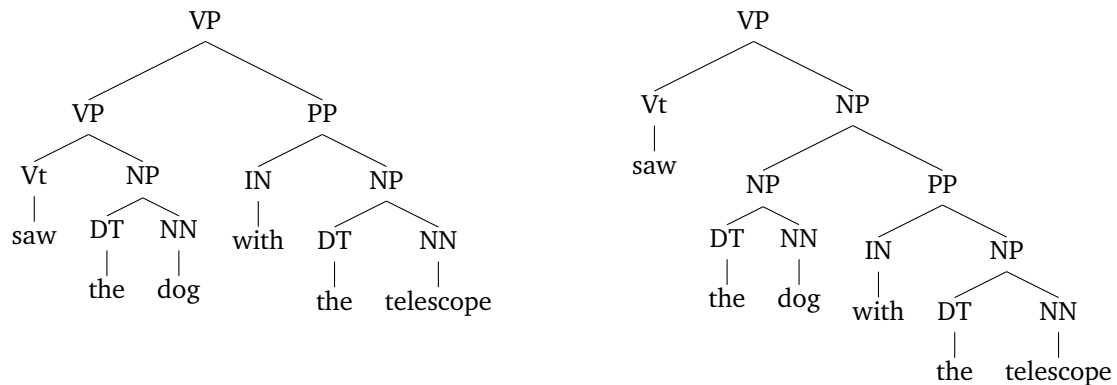


Sentence 3, Interpretation 2:
Flies from a place named Monty are fond of mosquitoes



Exercise 2.

Say we have the phrase *saw the dog with the telescope* and we are given the gold parse tree (left) and the predicted parse tree (right):



What are the precision and recall of this predicted parse tree?

SOLUTION

The gold tree has 11 nonterminal nodes, each expanded using certain rule.

The predicted tree has 11 nonterminal nodes, each expanded using a rule that may be or not the right one. In particular, the predicted rules $VP \rightarrow Vt NP$ and $NP \rightarrow NP PP$ are not in the gold tree. All the other rules are found in the gold tree, so they are right.

Precision is the number of right rules divided by the number of predicted rules (i.e. size of the predicted tree). So in this case, it is $9/11 = 81.8\%$.

Recall is the number of right rules divided by the number of expected rules (i.e. size of the gold tree). So in this case, it is also $9/11 = 81.8\%$ (since the size of both trees was the same, which shouldn't necessarily happen)

Exercise 3.

Consider the following CFG:

$S \rightarrow NP VP$	$DT \rightarrow \text{the}$	$NN \rightarrow \text{park}$
$NP \rightarrow DT NN$	$NN \rightarrow \text{man}$	$VB \rightarrow \text{saw}$
$NP \rightarrow NP PP$	$NN \rightarrow \text{dog}$	$IN \rightarrow \text{with}$
$PP \rightarrow IN NP$	$NN \rightarrow \text{cat}$	$IN \rightarrow \text{under}$
$VP \rightarrow VB NP$		

1. How many parse trees are there under this grammar for the sentence *the man saw the dog in the park* ?
2. How many parse trees are there under this grammar for the sentence *the man saw the dog in the park with the cat* ?

SOLUTION

1. The sentence *the man saw the dog in the park* has a unique analysis in this grammar, where the dog is in the park. This is because the grammar does not have a rule such as e.g. $VP \rightarrow VP PP$ that allows the PP *in the park* to be attached to the verb *saw*.
2. The sentence *the man saw the dog in the park with the cat* has two analysis under this grammar: One where the cat is with the dog, and another where the cat is with the park.

Exercise 4.

Consider the following CFG:

$S \rightarrow NP VP$	$DT \rightarrow the$	$NNS \rightarrow cats$
$NP \rightarrow DT NN$	$NN \rightarrow man$	$NNS \rightarrow parks$
$NP \rightarrow DT NNS$	$NN \rightarrow dog$	$VB \rightarrow see$
$NP \rightarrow NP PP$	$NN \rightarrow cat$	$VB \rightarrow sees$
$PP \rightarrow IN NP$	$NN \rightarrow park$	$IN \rightarrow in$
$VP \rightarrow VB NP$	$NNS \rightarrow dogs$	$IN \rightarrow with$
$VP \rightarrow VP PP$		

This grammar overgenerates incorrect English sentences, such as:

the dog see the cat
the dog in the park see the cat
the dog in the park see the cat in the park
the dogs sees the cat
the dogs in the park sees the cat
the dogs in the park sees the cat in the park

1. Modify the grammar so that all generated sentences respect third-person subject-verb agreement rules for English

SOLUTION

The rule joining the subject and the verb of the sentences is $S \rightarrow NP VP$, so we need to alter this rule to allow only the combination of singular NP with third person VP, and plural NP with non-third person VP. For this, we need different rules for singular/plural NP and for third/non-third person VP.

Thus, the top rule $S \rightarrow NP VP$ needs to be replaced with:

$S \rightarrow NP_s VP_s$
 $S \rightarrow NP_p VP_p$

All the NP rules must distinguish both kinds of noun phrases, replacing them with:

$NP_s \rightarrow DT NN$
 $NP_p \rightarrow DT NNS$
 $NP_s \rightarrow NP_s PP$
 $NP_p \rightarrow NP_p PP$

Finally, the rules for verb phrases must also distinguish both cases:

$VB_s \rightarrow sees$
 $VB_p \rightarrow see$
 $VP_s \rightarrow VB_s NP$
 $VP_p \rightarrow VB_p NP$
 $VP_s \rightarrow VP_s PP$
 $VP_p \rightarrow VP_p PP$

To avoid an explosion of rules, we can keep a generic NP to be used for noun phrases after the verb or inside a PP:

$NP \rightarrow NP_s$
 $NP \rightarrow NP_p$

Probabilistic Context Free Grammars

Exercise 5.

Using the following PCFG in CNF:

$S \rightarrow NP VP$	1.0	$P \rightarrow with$	1.0
$NP \rightarrow NP PP$	0.4	$V \rightarrow saw$	1.0
$PP \rightarrow P NP$	1.0	$NP \rightarrow astronomers$	0.1
$VP \rightarrow V NP$	0.7	$NP \rightarrow ears$	0.18
$VP \rightarrow VP PP$	0.3	$NP \rightarrow saw$	0.04
		$NP \rightarrow stars$	0.18
		$NP \rightarrow telescopes$	0.1

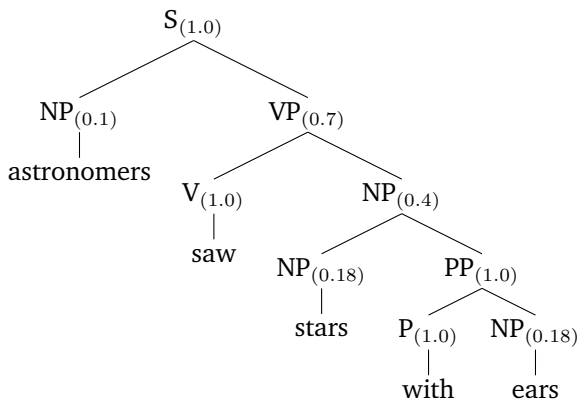
Work with the sentence: *astronomers saw stars with ears*

- How many correct parses are there for this sentence?
- Write them, along with their probabilities.

SOLUTION

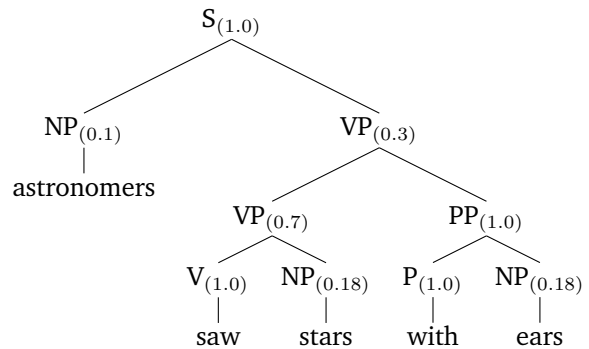
There are two possible parsers for this sentence according to the given grammar:

Option 1: (the stars had ears)



Probability: $1.0 \times 0.1 \times 0.7 \times 1.0 \times 0.4 \times 0.18 \times 1.0 \times 1.0 \times 0.18 = 0.00091$

Option 2: (Astronomers had their ears while watching the stars –or used ears to watch them)



Probability: $1.0 \times 0.1 \times 0.3 \times 0.7 \times 1.0 \times 0.18 \times 1.0 \times 1.0 \times 0.18 = 0.00068$

Exercise 6.

Given the following PCFG:

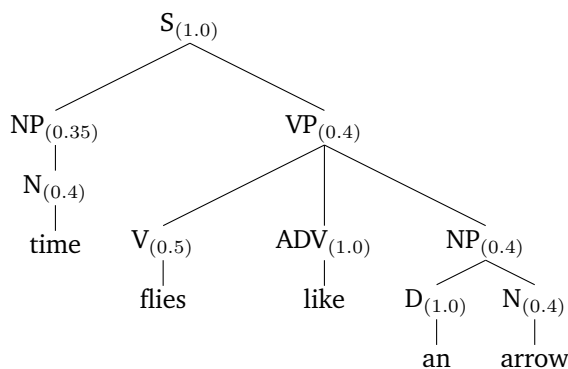
$S \rightarrow NP VP$	1.0	$N \rightarrow \text{time}$	0.4
$NP \rightarrow N N$	0.25	$N \rightarrow \text{flies}$	0.2
$NP \rightarrow D N$	0.4	$N \rightarrow \text{arrow}$	0.4
$NP \rightarrow N$	0.35	$D \rightarrow \text{an}$	1.0
$VP \rightarrow V NP$	0.6	$ADV \rightarrow \text{like}$	1.0
$VP \rightarrow V ADV NP$	0.4	$V \rightarrow \text{flies}$	0.5
		$V \rightarrow \text{like}$	0.5

and the sentence *time flies like an arrow*

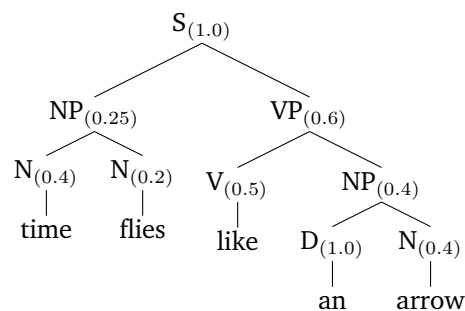
- Write two parse trees that this grammar generates for this sentence
- Compute the probability of each tree.
- Convert the grammar to CNF and emulate the behaviour of the CKY algorithm on this sentence. Provide the final chart with all the information involved.

SOLUTION

- Option 1: (time goes by so fast that reminds of an arrow)



- Option 2: (Alien 4-dimensional flies are fond of arrows)



- Option 1 probability: $1.0 \times 0.35 \times 0.4 \times 0.4 \times 0.5 \times 1.0 \times 0.4 \times 1.0 \times 0.4 = 0.00448$
Option 2 Probability: $1.0 \times 0.25 \times 0.4 \times 0.2 \times 0.6 \times 0.5 \times 0.4 \times 1.0 \times 0.4 = 0.00096$

- (a) Conversion of the grammar to CNF:

Chomsky Normal Form requires that all rules have a right hand side with exactly two non-terminals, or exactly one terminal. Rules $NP \rightarrow N$ and $VP \rightarrow V ADV NP$ violate this condition, so we need to transform them.

On the one hand, we shrink rule $NP \rightarrow N$ with those which produce N , so with $N \rightarrow \text{time}$, $N \rightarrow \text{flies}$ and $N \rightarrow \text{arrow}$ and finally we include the following rules:

$$NP \rightarrow \text{time} \quad 0.35 \times 0.4 = 0.14$$

$$NP \rightarrow \text{flies} \quad 0.35 \times 0.2 = 0.07$$

$$NP \rightarrow \text{arrow}, i \quad 0.35 \times 0.4 = 0.14$$

We remove the original rule not in CNF, but we keep the other 3 rules as they are used by other rules (ex: $NP \rightarrow N N$).

On the other hand, we split rule $VP \rightarrow V ADV NP$ to get two new rules in CNF:

$$VP \rightarrow V ADVP \quad 0.4$$

$$ADVP \rightarrow ADV NP \quad 1.0$$

The resulting grammar in CNF is:

$S \rightarrow NP VP$	1.0	$N \rightarrow time$	0.4
$NP \rightarrow N N$	0.25	$N \rightarrow flies$	0.2
$NP \rightarrow D N$	0.4	$N \rightarrow arrow$	0.4
$NP \rightarrow time$	0.14	$D \rightarrow an$	1.0
$NP \rightarrow flies$	0.07	$ADV \rightarrow like$	1.0
$NP \rightarrow arrow$	0.14	$V \rightarrow flies$	0.5
$VP \rightarrow V NP$	0.6	$V \rightarrow like$	0.5
$VP \rightarrow V ADVP$	0.4		
$ADVP \rightarrow ADV NP$	1.0		

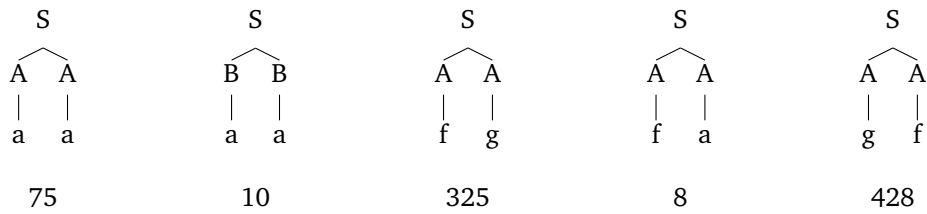
(b) CKY chart:

				15 0.00448 $S \rightarrow N_{11}VP_{25}$ ($1.0 \times 0.14 \times 0.032$) 0.00096 $S \rightarrow NP_{12}VP_{35}$ ($1.0 \times 0.02 \times 0.048$)
			14 0.0033 $S \rightarrow N_{11}VP_{25}$ ($1.0 \times 0.07 \times 0.048$) 0.032 $VP \rightarrow V_{22}ADVP_{35}$ ($0.4 \times 0.5 \times 0.16$)	25 0.0033 $S \rightarrow N_{11}VP_{25}$ ($1.0 \times 0.07 \times 0.048$) 0.032 $VP \rightarrow V_{22}ADVP_{35}$ ($0.4 \times 0.5 \times 0.16$)
		13	24 0.16 $ADVP \rightarrow ADV_{33}NP_{45}$ ($1.0 \times 1.0 \times 0.16$) 0.048 $VP \rightarrow V_{33}NP_{45}$ ($0.6 \times 0.5 \times 0.16$)	35 0.16 $ADVP \rightarrow ADV_{33}NP_{45}$ ($1.0 \times 1.0 \times 0.16$) 0.048 $VP \rightarrow V_{33}NP_{45}$ ($0.6 \times 0.5 \times 0.16$)
	12 0.02 $NP \rightarrow N_{11}N_{22}$ ($0.25 \times 0.4 \times 0.2$)	23	34 0.16 $NP \rightarrow D_{44}N_{55}$ ($0.4 \times 1.0 \times 0.4$)	45 0.16 $NP \rightarrow D_{44}N_{55}$ ($0.4 \times 1.0 \times 0.4$)
11 0.14 $NP \rightarrow time$ 0.4 $N \rightarrow time$	22 0.07 $NP \rightarrow flies$ 0.2 $N \rightarrow flies$ 0.5 $V \rightarrow flies$	33 1.0 $ADV \rightarrow like$ 0.5 $V \rightarrow like$	44 1.0 $D \rightarrow an$	55 0.4 $N \rightarrow arrow$
time	flies	like	an	arrow

(Blue rule in cell 35 indicates the most likely subtree selected in that cell)

Exercise 7.

Consider that you have as a training corpus a treebank containing the following trees. Each tree was observed the number of times indicated below it.



1. What PCFG would one get from this treebank (using MLE)?
2. Given the obtained grammar:
 - What is the most likely parse of the string aa ?
 - Is this a reasonable result? Discuss why.

SOLUTION

1. The given collection of training trees, taking into account the number of repetitions of each, will produce the following counts of rule applications:

$S \rightarrow AA$	$1 \times 75 + 0 \times 10 + 1 \times 325 + 1 \times 8 + 1 \times 428 = 836$
$S \rightarrow BB$	$0 \times 75 + 1 \times 10 + 0 \times 325 + 0 \times 8 + 0 \times 428 = 10$
$S \rightarrow \text{anything}$	$1 \times 75 + 1 \times 10 + 1 \times 325 + 1 \times 8 + 1 \times 428 = 846$
$B \rightarrow a$	$0 \times 75 + 2 \times 10 + 0 \times 325 + 0 \times 8 + 0 \times 428 = 20$
$B \rightarrow \text{anything}$	$0 \times 75 + 2 \times 10 + 0 \times 325 + 0 \times 8 + 0 \times 428 = 20$
$A \rightarrow a$	$2 \times 75 + 0 \times 10 + 0 \times 325 + 1 \times 8 + 0 \times 428 = 158$
$A \rightarrow f$	$0 \times 75 + 0 \times 10 + 1 \times 325 + 1 \times 8 + 1 \times 428 = 761$
$A \rightarrow g$	$0 \times 75 + 0 \times 10 + 1 \times 325 + 0 \times 8 + 1 \times 428 = 753$
$A \rightarrow \text{anything}$	$2 \times 75 + 0 \times 10 + 2 \times 325 + 2 \times 8 + 2 \times 428 = 1672$

Thus, the MLE probability for each rule would be:

$P(S \rightarrow AA) = P(AA S) = \#(S \rightarrow AA)/\#(S \rightarrow \text{anything}) = 836/846 = 0.988$
$P(S \rightarrow BB) = P(BB S) = \#(S \rightarrow BB)/\#(S \rightarrow \text{anything}) = 10/846 = 0.012$
$P(B \rightarrow a) = P(a B) = \#(B \rightarrow a)/\#(B \rightarrow \text{anything}) = 20/20 = 1.000$
$P(A \rightarrow a) = P(a A) = \#(A \rightarrow a)/\#(A \rightarrow \text{anything}) = 158/1672 = 0.095$
$P(A \rightarrow f) = P(f A) = \#(A \rightarrow f)/\#(A \rightarrow \text{anything}) = 761/1672 = 0.455$
$P(A \rightarrow g) = P(g A) = \#(A \rightarrow g)/\#(A \rightarrow \text{anything}) = 753/1672 = 0.450$

2. The input sequence aa can be derived by the obtained grammar in only two ways, which correspond to the first two trees in the training data.

The first tree has probability: $P(S \rightarrow AA) \times P(A \rightarrow a) \times P(A \rightarrow a) = 0.988 \times 0.095 \times 0.095 = 0.009$

The second tree has probability: $P(S \rightarrow BB) \times P(B \rightarrow a) \times P(B \rightarrow a) = 0.012 \times 1.000 \times 1.000 = 0.012$

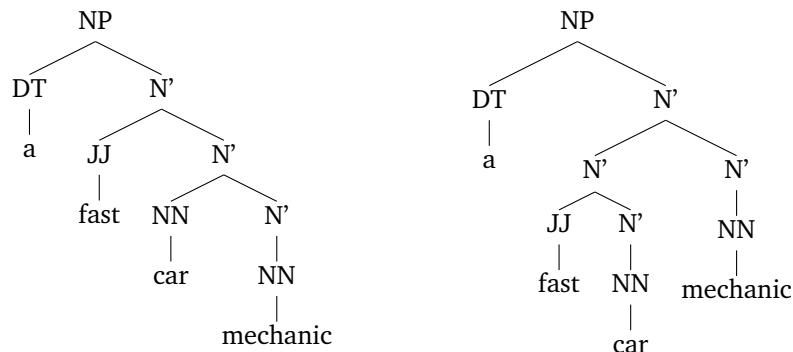
So, the most likely parse tree is the second one.

The first tree appears 75 times in the training data, while the second one occurs only 10 times, so one would expect the probability for the former to be higher. However, we do not compute the tree probability by counting how many times the whole tree occurs, but we just approximate it by multiplying the individual rule probabilities. Thus, the fact that B produces a with much higher probability than A is biasing the result.

This is due to the reduced amount of training data combined with the use of MLE: Since we do not perform any smoothing to consider the possibility of B producing other symbols, we are overestimating the probability of the rule $B \rightarrow a$.

Exercise 8.

Consider the two following parse trees:



Discuss whether the following statements are true or false and why:

1. The two parse trees receive the same probability under any PCFG
2. The first parse tree receives higher probability if $P(N' \rightarrow NN N') > P(N' \rightarrow N' N')$
3. The first parse tree receives higher probability if $P(N' \rightarrow NN N') > P(N' \rightarrow N' N') + P(N' \rightarrow NN)$

SOLUTION

1. False, since the set of rules used in each tree differ: left tree uses rule $N' \rightarrow NN N'$ while the tree on the right uses the rule $N' \rightarrow N' N'$. Also, the former uses rule $N' \rightarrow NN$ once but the latter uses it twice. So, assuming Q is the product of probabilities of rules shared by both trees, the first tree has probability $Q \times P(N' \rightarrow NN N')$, and the second has probability $Q \times P(N' \rightarrow N' N') \times P(N' \rightarrow NN)$, which are not necessarily equal.
2. True, since under this condition $Q \times P(N' \rightarrow NN N') > Q \times P(N' \rightarrow N' N')$. If we multiply the right hand side term by $P(N' \rightarrow NN)$ which is smaller than 1, the difference will increase.
3. False, since probabilities are multiplied, not added. The first tree would have higher probability if $P(N' \rightarrow NN N') > P(N' \rightarrow N' N') \times P(N' \rightarrow NN)$

Exercise 9.

Consider the following PCFG

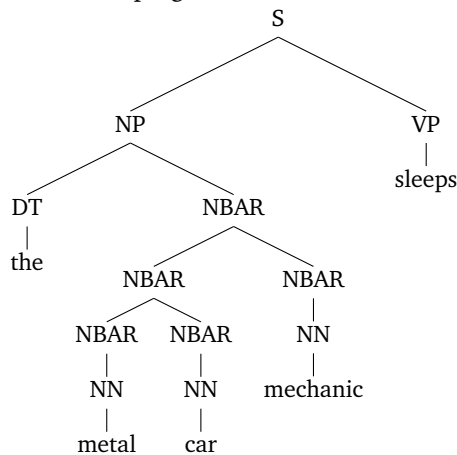
$S \rightarrow NP VP$	1.0	$VP \rightarrow sleeps$	1.0
$NP \rightarrow DT NBAR$	1.0	$DT \rightarrow the$	1.0
$NBAR \rightarrow NN$	0.7	$NN \rightarrow mechanic$	0.1
$NBAR \rightarrow NBAR NBAR$	0.3	$NN \rightarrow car$	0.2
		$NN \rightarrow metal$	0.7

1. What is the parse tree with highest probability for the sentence *the metal car mechanic sleeps* ?
2. Modify the grammar above so that the sentence *the human language technology rules* has two interpretations (one about *human language* and another about *human technology*). Draw the trees for both interpretations, and point out which is the most likely.

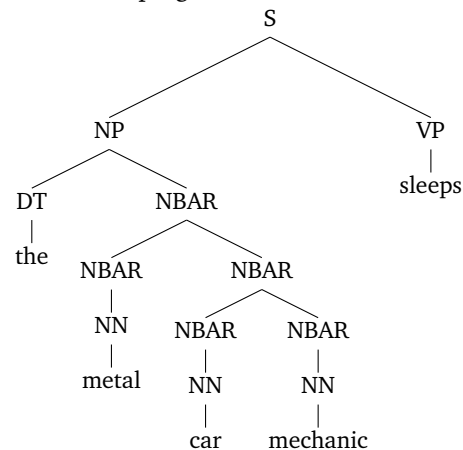
SOLUTION

1. This grammar produces two possible trees for the given sentence:

Option 1: (the mechanic that works on metal cars is sleeping)



Option 2: (the metal mechanic that works on cars is sleeping)



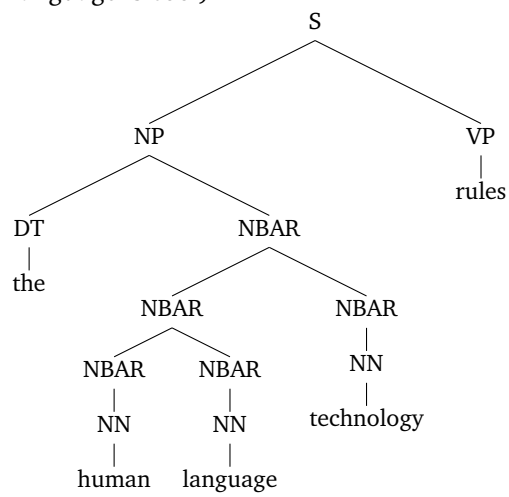
Both trees use once the rules $S \rightarrow NP VP$ and $NP \rightarrow DT NBAR$, twice the rule $NBAR \rightarrow NBAR NBAR$, and three times the rule $NBAR \rightarrow NN$. Also, the rules producing the leaves are also the same. The only difference is the order in which the rules are applied. Thus, the probabilities of both trees are identical, and there is not one single best tree, but two.

2. The grammar already allows the ambiguous structure for the sequence NN NN NN. We only need to add the new words to the grammar, and fix the probabilities. Rule probabilities are invented, since we do not have training data, but we need to ensure that the rules for the same non-terminal symbol add up to 1. New rules are highlighted in red. Redistributed probabilities are shown in blue.

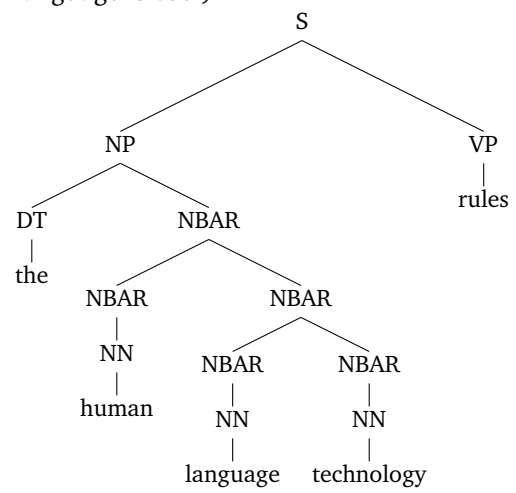
$S \rightarrow NP VP$	1.0	$VP \rightarrow sleeps$	0.5	$VP \rightarrow rules$	0.5
$NP \rightarrow DT NBAR$	1.0	$DT \rightarrow the$	1.0		
$NBAR \rightarrow NN$	0.7	$NN \rightarrow mechanic$	0.1	$NN \rightarrow language$	0.2
$NBAR \rightarrow NBAR NBAR$	0.3	$NN \rightarrow car$	0.2	$NN \rightarrow technology$	0.2
		$NN \rightarrow metal$	0.2	$NN \rightarrow human$	0.1

Possible trees for the new sentence are :

Option 1: (technology that deals with human language is cool)



Option 2: (human technology that deals with language is cool)



Again, since both trees have the same rules, they have exactly the same probability.