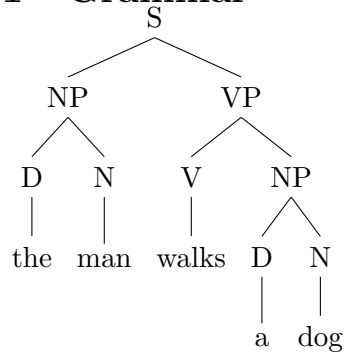


NLP: Parsing

Dan Garrette
dhg@cs.utexas.edu

December 27, 2013

1 Grammar



- “Constituency” parse
- S (sentence), NP (noun phrase), VP (verb phrase) are constituents
- Words combine to make phrases, and phrases combine to make larger phrases and sentences.

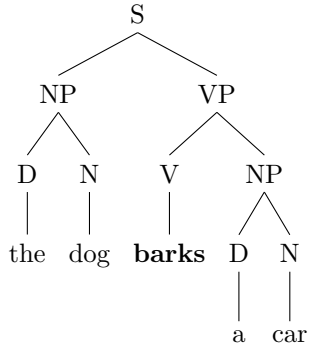
2 Context-Free Grammars (CFG)

- Context-free grammars can be specified by a table of “productions”

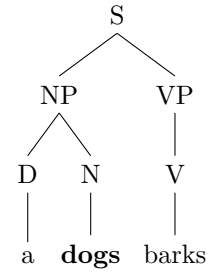
S	→	NP VP		D	→	{ <i>the, a</i> }
NP	→	D N	(“ <u>a dog</u> barks”)	N	→	{ <i>man, dog</i> }
NP	→	N	(“ <u>dogs</u> bark”)			
VP	→	V NP	(transitive verb)	V	→	{ <i>barks, walks, saw</i> }
VP	→	V	(intransitive verb)			

- Words are called “terminals” and other nodes are “non-terminals”

- Independence assumption: each rule choice is dependent only on the parent node; it is independent of all other rule choices
- Independence assumption leads to many problems. A couple:



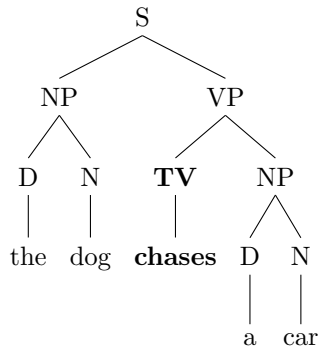
(a) Intransitive verb with a direct object



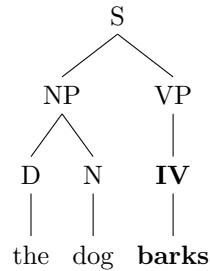
(b) Noun doesn't agree in number with determiner or verb

- We can solve some of these issues by refining our CFG. For example:

VP	→	IV	(intransitive verb)	IV	→	{ <i>barks, walks</i> }
VP	→	TV NP	(transitive verb)	TV	→	{ <i>chases, walks</i> }



(c) Intransitive verbs are disallowed



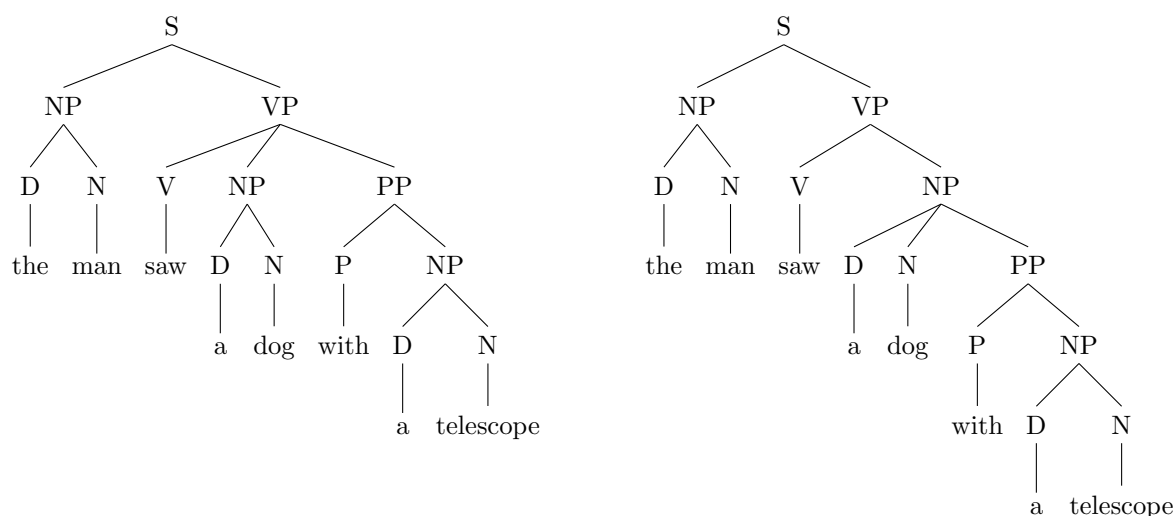
(d) Transitive verbs are disallowed

3 Syntactic Ambiguity

- For a given CFG, there can be multiple trees that describe the same sentence
- Add the following rules to the above:

$NP \rightarrow D N PP$ (prep phrase attach to noun) $N \rightarrow \{telescope\}$
 $VP \rightarrow V NP PP$ (prep phrase attach to verb)
 $PP \rightarrow P NP$ ("in the house")

- "The man saw a dog with a telescope"



- Ambiguity is explosive
- A sentence ending in n prepositional phrases has *over* 2^n parses (catalan numbers)
 - "I saw the man with the telescope": 2 parses
 - "I saw the man on the hill with the telescope": 5 parses
 - "I saw the man on the hill in texas with the telescope": 14 parses
 - "I saw the man on the hill in texas with the telescope at noon": 42 parses
 - "I saw the mon on the hill in texas with the telescope at noon on monday": 132 parses

4 Agreement

- In order for a sentence to be grammatical, we must also respect *agreement rules*
 - number: *a dog* vs. *all dogs*
 - person: 1st person (*I am*), 2nd person (*you are*), 3rd person (*he is*)
 - gender: *un homme* vs. *une femme* or even *she sees herself*

- The grammar defined above does not enforce this
 - For an NP, since productions are independent of each other, we could get either *dog* or *dogs* no matter whether the D is *a* or *all*
- We can incorporate this into our grammar by duplicating our productions to ensure agreement:

S	→	NP _{sg} VP _{sg}	D _{sg}	→	{ <i>the, a</i> }
S	→	NP _{pl} VP _{pl}	D _{pl}	→	{ <i>the, all</i> }
NP _{pl}	→	N _{pl} (“ <u>dogs</u> bark”)	N _{sg}	→	{ <i>dog</i> }
NP _{pl}	→	D _{pl} N _{pl} (“ <u>all dogs</u> bark”)	N _{pl}	→	{ <i>dogs</i> }
NP _{sg}	→	D _{sg} N _{sg} (“ <u>a dog</u> barks”)			
VP _{sg}	→	V _{sg} (“a dog <u>barks</u> ”)	V _{sg}	→	{ <i>barks, walks, sees</i> }
VP _{sg}	→	V _{sg} NP _{sg} (“a dog <u>chases a car</u> ”)	V _{pl}	→	{ <i>bark, walk, see</i> }
VP _{sg}	→	V _{sg} NP _{pl} (“a dog <u>chases cars</u> ”)			
VP _{pl}	→	V _{pl} (“all dogs <u>bark</u> ”)			
VP _{pl}	→	V _{pl} NP _{sg} (“all dogs <u>chase a car</u> ”)			
VP _{pl}	→	V _{pl} NP _{pl} (“all dogs <u>chase cars</u> ”)			

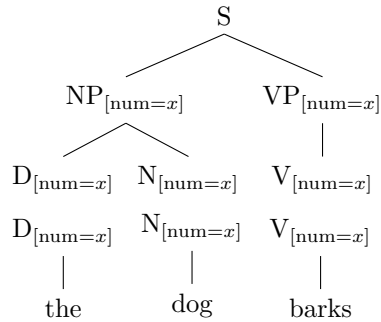
- But this quickly explodes the number of rules
- And the explosion is worse with more agreement rules:

VP _{sg,1st,masc}	→	V _{sg,1st,masc}
VP _{sg,1st,fem}	→	V _{sg,1st,fem}
VP _{sg,2nd,masc}	→	V _{sg,2nd,masc}
VP _{sg,2nd,fem}	→	V _{sg,2nd,fem}
...		
VP _{pl,3rd,fem}	→	V _{pl,3rd,fem}

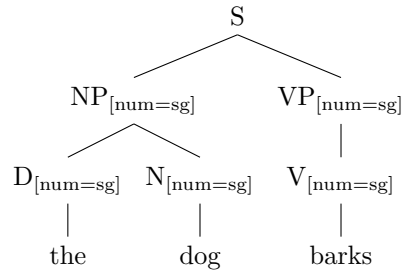
- We can simplify this greatly using *feature structures* that use variables to ensure agreement

S	→	NP _[num=x] VP _[num=x]	D _[num=sg]	→	{ <i>a, every</i> }
			D _[num=pl]	→	{ <i>all, some</i> }
NP _[num=x]	→	D _[num=x] N _[num=x]	D _[num=z]	→	{ <i>the</i> }
VP _[num=x]	→	V _[num=x]	N _[num=sg]	→	{ <i>dog</i> }
VP _[num=x]	→	V _[num=x] NP _[num=y]	N _[num=pl]	→	{ <i>dogs</i> }
			V _[num=sg]	→	{ <i>barks, walks, sees</i> }
			V _[num=pl]	→	{ <i>bark, walk, see</i> }

- Feature structures must *unify* as they are combined
 - For example, all x s and z must resolve to the same value (sg or pl)



(e) Need to unify underspecified features in the rules with specified features in the N and V terminal rules



(f) all x s and z resolve to “sg”

- Similar features could be set up for other required agreements
 - $\text{PRO}[\text{num}=\text{sg}, \text{per}=\text{3rd}, \text{gen}=\text{masc}] \rightarrow \text{it}$

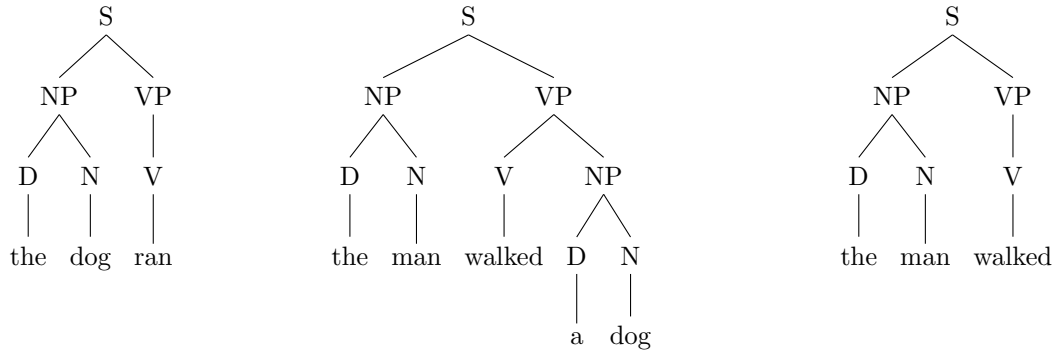
5 Probabilistic Context-Free Grammars (PCFG)

- A CFG is deterministic
 - It can decide whether a sentence is in the language (grammatical), or not
 - It can't judge whether one sentence is more likely than another
 - Problematic since we want to say that every sentence is *possible*, even if it's not likely
- A PCFG is a CFG in which, for every non-terminal, we have a probability distribution over possible productions
- In other words, for each non-terminal A , we have a distribution $p(\beta \mid A)$
 - The probability that, given a parent non-terminal A , we choose the production rule that yields β
 - We can alternatively write this as $p(A \rightarrow \beta)$

S	→	NP VP	1.0	D	→	<i>the</i>	0.6
				D	→	<i>a</i>	0.4
NP	→	D N	0.7	N	→	<i>man</i>	0.5
NP	→	N	0.2	N	→	<i>dog</i>	0.4
NP	→	D N PP	0.1	N	→	<i>telescope</i>	0.1
VP	→	V NP	0.4	V	→	<i>barks</i>	0.2
VP	→	V	0.4	V	→	<i>walks</i>	0.4
VP	→	V NP PP	0.2	V	→	<i>saw</i>	0.4

Estimating parameters (MLE)

- We can calculate the MLE parameters of a PCFG model by counting productions in a corpus of parse trees
- Assume these three sentences comprise a corpus:



- We estimate by counting up all productions in the corpus and normalizing
 - Since every non-terminal A must produce *something*, we can simplify slightly

$$p(\beta \mid A) = \frac{C(A \rightarrow \beta)}{\sum_{\beta' \in P} C(A \rightarrow \beta')} = \frac{C(A \rightarrow \beta)}{C(A)}$$

- Estimates from the above corpus of trees yield:

			$C(A \rightarrow \beta)$	$p(\beta \mid A)$			$C(A \rightarrow \beta)$	$p(\beta \mid A)$
S	→	NP VP	3	1.0	D	→	<i>the</i>	0.75
					D	→	<i>a</i>	0.25
NP	→	D N	4	1.0	N	→	<i>dog</i>	0.5
VP	→	V	2	0.66	N	→	<i>man</i>	0.5
VP	→	V NP	1	0.33	V	→	<i>walked</i>	0.66
					V	→	<i>ran</i>	0.33

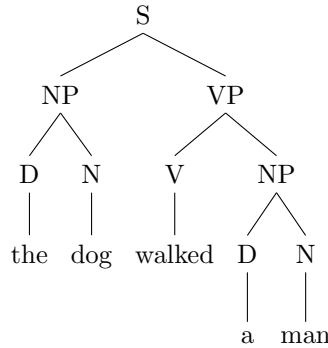
Add- λ smoothing

- Same as always (where P is the set of all possible β s produced):

$$p(\beta \mid A) = \frac{C(A \rightarrow \beta) + \lambda}{\sum_{\beta' \in P} (C(A \rightarrow \beta') + \lambda)} = \frac{C(A \rightarrow \beta) + \lambda}{(\sum_{\beta' \in P} C(A \rightarrow \beta')) + \lambda|P|} = \frac{C(A \rightarrow \beta) + \lambda}{C(A) + \lambda|P|}$$

Likelihood of a (parsed) sentence

- We can use this to calculate the likelihood of seeing a particular parse of a particular sentence
- Multiply the probabilities of all productions found in the tree
- Assume this tree:



- Count up the number of each production in the tree, and get the probability of each production from the above table

			count	prob				count	prob
S	→	NP VP	1	1.0	D	→	<i>the</i>	1	0.75
					D	→	<i>a</i>	1	0.25
NP	→	D N	2	1.0	N	→	<i>dog</i>	1	0.5
VP	→	V NP	1	0.33	N	→	<i>man</i>	1	0.5
					V	→	<i>walked</i>	1	0.66

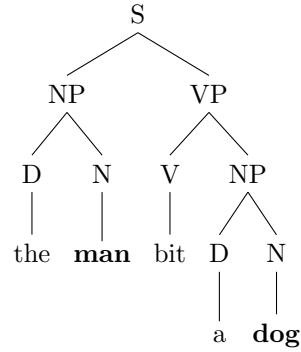
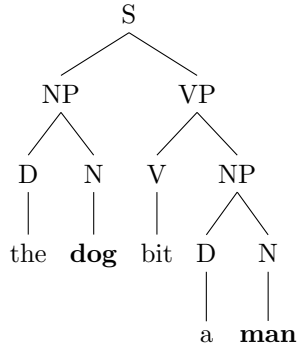
- Find the product:

$$\begin{aligned}
& p(S \rightarrow NP VP) \cdot p(NP \rightarrow D N)^2 \cdot p(VP \rightarrow V NP) \cdot p(D \rightarrow the) \cdot p(D \rightarrow a) \cdot p(N \rightarrow dog) \\
& \quad \cdot p(N \rightarrow man) \cdot p(V \rightarrow walked) \\
& = 1.0 \cdot 1.0^2 \cdot 0.33 \cdot 0.75 \cdot 0.25 \cdot 0.5 \cdot 0.5 \cdot 0.66
\end{aligned}$$

- Note that the $p(NP \rightarrow D N)$ term is squared since the production “ $NP \rightarrow D N$ ” appears twice in the tree

6 Lexicalized Grammars

- Problem: trees do not take semantic coherence into account
- Production rules are independent



- These two sentence have *exactly* the same likelihood
- Since it's a product, we can swap two productions without changing the result

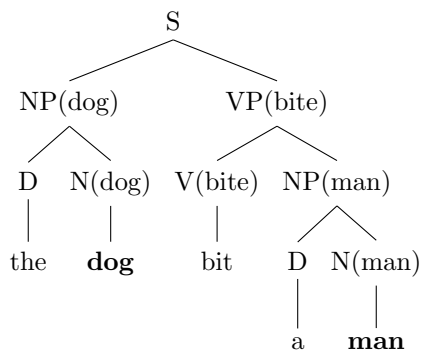
$$\begin{aligned}
 & p(S \rightarrow NP VP) \cdot p(NP \rightarrow D N)^2 \cdot p(VP \rightarrow V NP) \cdot p(D \rightarrow the) \cdot p(D \rightarrow a) \cdot p(N \rightarrow \textit{dog}) \cdot p(N \rightarrow \textit{man}) \cdot p(V \rightarrow bit) \\
 & \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \downarrow \qquad \qquad \qquad \downarrow \\
 & p(S \rightarrow NP VP) \cdot p(NP \rightarrow D N)^2 \cdot p(VP \rightarrow V NP) \cdot p(D \rightarrow the) \cdot p(D \rightarrow a) \cdot p(N \rightarrow \textit{man}) \cdot p(N \rightarrow \textit{dog}) \cdot p(V \rightarrow bit)
 \end{aligned}$$

- Solution: lexicalize the grammar
- Subcategorize rules with their head words:

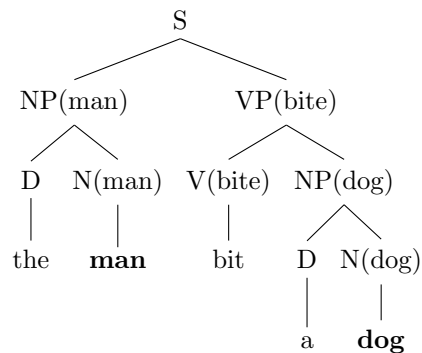
S	→	NP(man) VP(bite)	0.3	D	→	<i>the</i>	0.6
S	→	NP(dog) VP(bite)	0.7	D	→	<i>a</i>	0.4
S ¹	→	NP(sandwich) VP(bite)	0.0				
NP(man) ²	→	D N(man)	0.8	N(man) ²	→	<i>man</i>	0.7
NP(man) ²	→	N(man)	0.2	N(man) ²	→	<i>men</i>	0.3
NP(dog) ³	→	D N(dog)	0.4	N(dog) ³	→	<i>dog</i>	0.2
NP(dog) ³	→	N(dog)	0.6	N(dog) ³	→	<i>dogs</i>	0.8
VP(bite) ⁴	→	V(bite) NP(man)	0.2	V(bite)	→	<i>bite</i>	0.3
VP(bite) ⁴	→	V(bite) NP(dog)	0.1	V(bite)	→	<i>bites</i>	0.2
VP(bite) ⁴	→	V(bite) NP(sandwich)	0.7	V(bite)	→	<i>biting</i>	0.1
VP(bite) ⁵	→	V(bite)	0.0	V(bite)	→	<i>bit</i>	0.4

- ¹Sandwiches don't bite
- ²Men are likely talked about in the singular
- ³Dogs are likely talked about in the plural
- ⁴Men are bitten infrequently, dogs are bitten less frequently, and sandwiches are bitten often.
- ⁵"Biting" can't be intransitive

- Downside: increased sparsity!



(g) more likely



(h) unlikely

7 Generative Model

- Like naïve Bayes models, N-Gram models, and hidden Markov models
- Two probability distribution: $p(\beta \mid \alpha)$, for production rules $\alpha \rightarrow \beta$, and $p(\sigma)$, where σ is a possible “start” symbol
- Generative story:
 1. Choose a start symbol x from the distribution over start symbols $p(\sigma)$
 2. If x is a terminal, STOP
 3. Else, choose some β from $p(\beta \mid x)$
 4. For each symbol y in β , go to step 2
- For each node with symbol x , we choose a production rule of the form $x \rightarrow \beta$ according to their probabilities and then recursively choose rules for every node in β until we reach terminals for all branches.

8 Parsing

Top-Down

- Never try to build tree that doesn’t make a sentence.
- Waste time exploring trees that don’t end up with the correct words.

Bottom-Up

- Never make a tree that doesn’t start with the right words.
- Waste time exploring trees that don’t make a sentence.

9 Parsing with P-CKY

Find the mostly likely parse tree t for the sentence $w_0...w_n$

- $\hat{t} = \operatorname{argmax}_t p(t \mid w_0...w_n)$
- Since we know how to find the likelihood of a parse (from above), we could enumerate all possible trees, and find the likelihood of each one.
- But the number of possible trees is enormous
- We will use a *dynamic programming* algorithm to find the best parse tree: P-CKY
- Probabilistic version of the CKY algorithm, which can be used with any CFG
- Analogous to the Viterbi algorithm (both are dynamic programming algorithms)

Chomsky Normal Form (CNF)

- CKY requires that the grammar be in CNF
- All productions must be one of:
 - Non-terminal producing exactly 2 non-terminals: $A \rightarrow B C$
 - Non-terminal producing exactly 1 non-terminal: $A \rightarrow B$
 - Non-terminal producing exactly 1 terminal: $A \rightarrow w$
- Any CFG can be converted to CNF
 - Producing more than two nonterminals: “ $A \rightarrow B \dots Y Z$ ” becomes
 - * $A \rightarrow B \dots [Y+Z]$
 - $[Y+Z] \rightarrow Y Z$
 - * Repeat as necessary until you just have $A \rightarrow B X$
 - * $NP \rightarrow D \text{ Adj } N$ becomes
 - $NP \rightarrow D [\text{Adj}+N]$
 - $[\text{Adj}+N] \rightarrow \text{Adj } N$
 - For our trees we should never have terminals and nonterminals mixed in productions since all terminals are produced by unary rules from POS tags, which never produce anything other than terminals. (But if we needed to we could create dummy non-terminals for those terminals.)

P-CKY Idea

1. Make a pass over the whole sentence, trying to find subtrees that explain spans within the sentence
2. Dynamic programming: store intermediate results in a chart
3. Start with the leftmost word. For each non-terminal (NT), find the most plausible subtree that covers just that word.

4. Move to the next word
 - a. Again, for each NT, find the most plausible subtree that covers just that word.
 - b. Then, for each NT, find the most plausible subtree that covers both the word and the previous word.
 - c. Same for the three-word span of the word and the previous two words.
 - d. Continue until you have the entire span from the beginning of the sentence to the word.
 - e. Then move to the next word and repeat.
5. You are finished when you have, for each NT, the most plausible tree covering the entire span of the sentence.
6. Pick the root node that is most plausible considering the *a priori* probability of an NT being the root of a sentence.

P-CKY Algorithm

- NT is the set of all non-terminals (including composites and POS tags)
- N = length of sentence $w_0...w_{N-1}$
- `table[i,j][A]` is the best possible score for a subtree spanning words $w_i...w_{j-1}$ with A as its root
- `back[i,j][A] = (k,B,C)` indicates that the best possible subtree spanning words $w_i...w_{j-1}$ with A as its root is comprised of a left subtree spanning words $w_i...w_{k-1}$ with B as its root and a right subtree spanning words $w_k...w_{j-1}$ with C as its root.

10 P-CKY Example

“the complex houses married students”

NP[the complex houses] VP[V[married] NP[N[students]]]
 NP[this complex] VP[V[houses] [NP[A[married] [N students]]]]

11 N-Best Parses

Useful for providing multiple possibilities to down-stream applications.

Can be fed into a discriminative re-ranker that uses something like MaxEnt with features built from linguistic knowledge to re-score the parses. The discriminative model wouldn't be used for *parsing*, just to calculate the likelihoods of the parses that were found.

12 Other Grammatical Formalisms

12.1 TAG: Tree-Adjoining Grammar

- <http://www.seas.upenn.edu/~joshi/joshi-schabes-tag-97.pdf>
- Instead of single-layer production rules of a CFG, production rules are tree *fragments*
- Fragments have “incomplete” nodes that must produce to other tree fragments. (Just like how “incomplete” nodes in a CFG must produce the next rule.)
- Lexicalized: Each fragment centered around a word
- Mildly context-sensitive (as opposed to context-free)

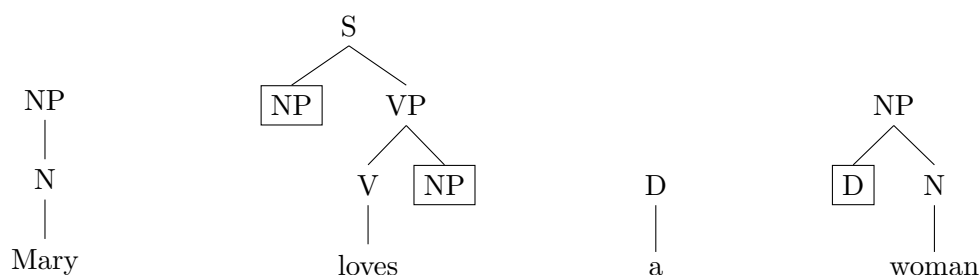


Figure 1: TAG productions

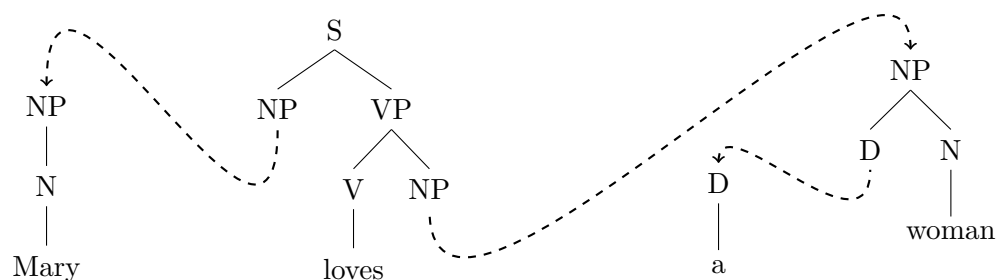


Figure 2: TAG fragment combination giving parse for “Mary loves a woman”

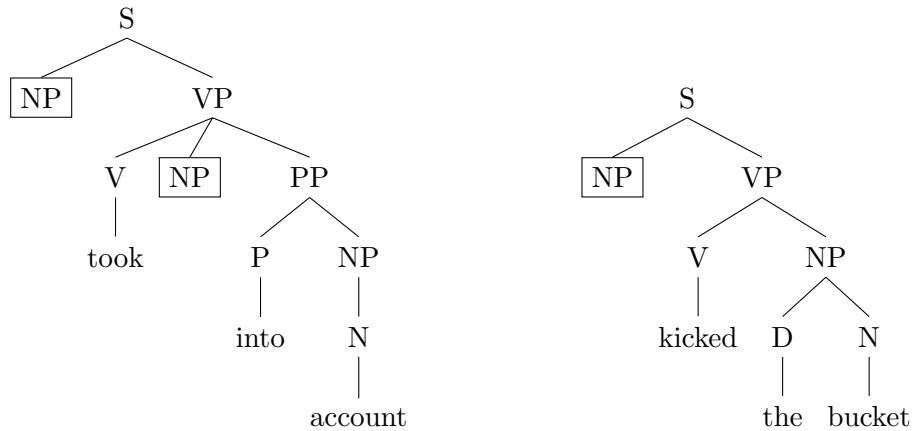


Figure 3: TAG productions for common phrases

12.2 CCG: Combinatory Categorical Grammar

- Instead of atomic POS tags and Non-Terminals, use *categories* that are constructed from other categories:
 - Categories are defined by a (context-free) grammar
 - atomic categories: $C \rightarrow \{ S, NP, N \}$
 - complex categories: $C \rightarrow \{ C/C, C \backslash C \}$
 - Slash and backslash operators indicate that the category is a function:
 - * A/B is a category that looks directly to its **right** for something of category B , and combines with it to produce an A
 - * $A \backslash B$ is a category that looks directly to its **left** for something of category B , and combines with it to produce an A
 - * So B is the *input* to the function, and A is the output
- Mildly context-sensitive
- Weakly equivalent to TAG

Words are assigned categories in the *Lexicon*:

Mary: NP

sleeps: $S \backslash NP$

loves: $(S \backslash NP) / NP$

a: NP / N

woman: N

The parse tree for a sentence, then, is the result of doing the combinations:

12.3 Dependency Parsing



Figure 4: Series of CCG combinations to parse “Mary sleeps”

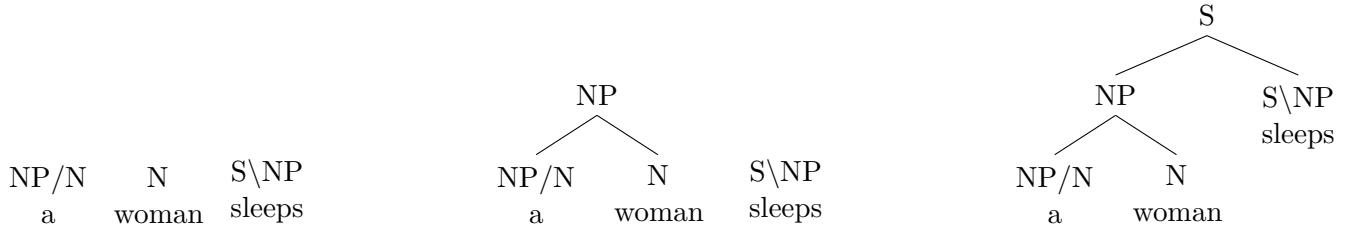


Figure 5: Series of CCG combinations to parse “a woman sleeps”

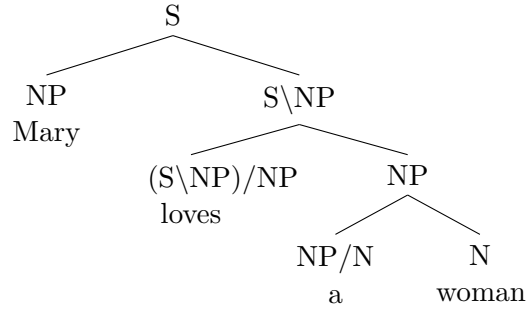


Figure 6: CCG parse tree for “Mary loves a woman”

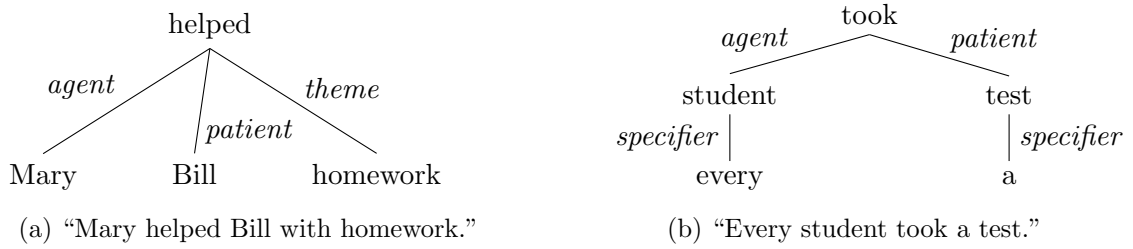


Figure 7: Dependency parse trees

Algorithm 1: P-CKY

Initialize **table** to be an $N+1 \times N+1$ array (for (best-so-far) scores)

Initialize **back** to be an $N+1 \times N+1$ array (for backpointers)

for $j \leftarrow 1$ **to** N **do**

// Fill in the bottom of the chart (single-word spans):

$i = j - 1$;

for $A \leftarrow NT$ (or restrict to just POS tags if convenient) **do**

// A is a potential POS tag for w_i

if $p(w_i | A) > 0.0$ **then**

// $A \rightarrow w_i$ is a valid rule

$\text{table}[i, j][a] = p(w_i | A)$

$\text{back}[i, j][a] = w_i$

// Fill in the higher levels of the chart (multi-word spans)

for $i \leftarrow j-1$ **downto** 0 **do**

// Binary Rules (with k as pivot)

for $k \leftarrow i+1$ **to** $j-1$ **do**

for $B \leftarrow \text{table}[i, k]$ and $C \leftarrow \text{table}[k, j]$ **do**

// B is the root of a sub-tree covering the span $w_i \dots w_{k-1}$

// C is the root of a sub-tree covering the span $w_k \dots w_{j-1}$

for $A \leftarrow NT$ (or restrict to just non-POS tags if convenient) **do**

// A is a potential root for the subtree spanning $w_i \dots w_{j-1}$

if $p(B \ C | A) > 0.0$ **then**

// $A \rightarrow B \ C$ is a valid rule

$s = p(B \ C | A) \cdot \text{table}[i, k][B] \cdot \text{table}[k, j][C]$

if $\text{table}[i, j]$ doesn't contain A **OR** $\text{table}[i, j][A] < s$ **then**

// Either we haven't seen a valid subtree covering i to j with root A ,

// or this new subtree has a higher score

$\text{table}[i, j][A] = s$

$\text{back}[i, j][A] = (k, B, C)$

// Unary Rules

 done = false

while !done **do**

 done = true

for $B \leftarrow \text{table}[i, j]$ **do**

// B is the root of a sub-tree covering the span $w_i \dots w_{j-1}$

for $A \leftarrow NT$ (or restrict to just non-POS, non-compound tags if convenient) **do**

// A is a potential root for the subtree spanning $w_i \dots w_{j-1}$

if $p(B | A) > 0.0$ **then**

// $A \rightarrow B$ is a valid rule

$s = p(B | A) \cdot \text{table}[i, j][B]$

if $\text{table}[i, j]$ doesn't contain A **OR** $\text{table}[i, j][A] < s$ **then**

// Either we haven't seen a valid subtree covering i to j with root A ,

// or this new subtree has a higher score

$\text{table}[i, j][A] = s$

$\text{back}[i, j][A] = B$

 done = false *// if we added a new NT, then re-check*

Algorithm 2: Retrieve the best-parse tree from the backpointer table

// Select the best root element while considering the prior over all potential root elements

root = $\underset{S \leftarrow \text{table}[0, N]}{\text{argmax}} \text{table}[0, N][S] \cdot p(S)$

if $\text{table}[0, N][\text{root}] \cdot p(\text{root}) = 0.0$ **then**

| *// There is no valid parse of the sentence*

else

| FollowBackpointers(root, 0, N)

Function FollowBackpointers(A, i, j) : Tree

back[i, j][A] **match**

case (k, B, C) \Rightarrow *// Binary branch*

 Tree(A, FollowBackpointers(B, i, k), FollowBackpointers(C, k, j))

case B \Rightarrow *// Unary branch*

 Tree(A, FollowBackpointers(B, i, j))

case w \Rightarrow *// Terminal (word)*

 Leaf(w)
