Probabilistic Context-Free Grammars Informatics 2A: Lecture 22

Adam Lopez

8 November, 2015

- CYK (non-probabilistic version)
- The Earley algorithm (three important operators that are run until we find a parse: PREDICTOR, SCANNER and COMPLETER)
- The complexity of CYK and the Earley algorithm: O(n³).
 (But, grammar constants can be large!)



2 Probabilistic Context-Free Grammars

- Definition
- Conditional Probabilities
- Applications
- Probabilistic CYK

Three things motivate the use of probabilities in grammars and parsing:

- Syntactic disambiguation. Ambiguity is unavoidable in natural language (grammars).
- Overage. What if there is a production that we haven't seen before? Might want to allow with small probability (Seemingly extreme but actually common use case: all productions are possible!).
- Representativeness. Suppose we want to adapt a parser to a new domain, where some words have different usage, hence different part-of-speech.

Motivation 1: Ambiguity

- Amount of ambiguity increases with sentence length.
- Real sentences are fairly long (avg. sentence length in the *Wall Street Journal* is 25 words).
- The amount of (unexpected!) ambiguity increases rapidly with sentence length. This poses a problem, even for chart parsers, if they have to keep track of all possible analyses.
- It would reduce the amount of work required if we could ignore improbable analyses.

A second provision passed by the Senate and House would eliminate a rule allowing companies that post losses resulting from LBO debt to receive refunds of taxes paid over the previous three years. [wsj_1822] (33 words)

- It is actually very difficult to write a grammar that covers all the constructions used in ordinary text or speech.
- Typically hundreds of rules are required in order to capture both all the different linguistic patterns and all the different possible analyses of the same pattern. (How many grammar rules did we have to add to cover three different analyses of You made her duck?)
- Ideally, one wants to induce (learn) a grammar from a corpus.
- Grammar induction requires probabilities.

The likelihood of a particular construction can vary, depending on:

- register (formal vs. informal): eg, greenish, alot, subject-drop (*Want a beer*?) are all more probable in informal than formal register;
- genre (newspapers, essays, mystery stories, jokes, ads, Twitter, etc.): PoS-taggers trained on different genres often have very different distributions.
- domain (biology, patent law, football, etc.).

Probabilistic grammars and parsers can reflect these kinds of distributions.

Book the dinner flight.











A PCFG $\langle N, \Sigma, R, S \rangle$ is defined as follows:

- N is the set of non-terminal symbols
- Σ is the terminals (disjoint from N)
- *R* is a set of rules of the form $A \rightarrow \beta[p]$ where $A \in N$ and $\beta \in (\sigma \cup N)$ *, and *p* is a number between 0 and 1
- S a start symbol, $S \in N$

A PCFG is a CFG in which each rule is associated with a probability.

What does the *p* associated with each rule express?

It expresses the probability that the LHS non-terminal will be expanded as the RHS sequence.

• $P(A \rightarrow \beta | A)$

•
$$\sum_{\beta} P(A \rightarrow \beta | A) = 1$$

• The sum of the probabilities associated with all of the rules expanding the non-terminal A is required to be 1.

$$A \rightarrow \beta \ [p]$$
 or $P(A \rightarrow \beta | A) = p$ or $P(A \rightarrow \beta) = p$

Example Grammar

$S \rightarrow NP VP$	[.80]	Det $ ightarrow$ the	[.10]
$S \rightarrow Aux NP VP$	[.15]	$\mathit{Det} ightarrow \mathit{a}$	[.90]
$S \rightarrow VP$	[.05]	Noun $ ightarrow$ book	[.10]
NP ightarrow Pronoun	[.35]	Noun $ ightarrow$ flight	[.30]
NP ightarrow Proper-Noun	[.30]	Noun $ ightarrow$ dinner	[.60]
NP ightarrow Det Nominal	[.20]	Proper-Noun $ ightarrow$ Houston	[.60]
NP ightarrow Nominal	[.15]	Proper-Noun $ ightarrow$ NWA	[.40]
Nominal $ ightarrow$ Noun	[.75]	Aux ightarrow does	[.60]
Nominal $ ightarrow$ Nominal Noun	[.05]	$\mathit{Aux} ightarrow \mathit{can}$	[.40]
VP ightarrow Verb	[.35]	$\mathit{Verb} ightarrow \mathit{book}$	[.30]
$VP ightarrow Verb \ NP$	[.20]	$\mathit{Verb} ightarrow \mathit{include}$	[.30]
$VP ightarrow Verb \ NP \ PP$	[.10]	Verb $ ightarrow$ prefer	[.20]
$VP \rightarrow Verb \ PP$	[.15]	$\mathit{Verb} ightarrow \mathit{sleep}$	[.20]

Start with the root node, and at each step, probabilistically expand the nodes until you hit a terminal symbol:

Qustion: Does this process always have to terminate?

Qustion: Does this process always have to terminate?

Consider the grammar, for some $\epsilon > 0$:

Example

- $S
 ightarrow S \, S$ with probability 0.5 + ϵ
- S
 ightarrow a with probability 0.5 ϵ

Qustion: Does this process always have to terminate?

Consider the grammar, for some $\epsilon > 0$:

Example	
${\it S} ightarrow {\it S} {\it S}$ with probability 0.5 $+ \epsilon$	
$S ightarrow a$ with probability 0.5 $- \epsilon$	

Whenever a nonterminal is expanded, it is more probable that the result will **increase** rather than **decrease** the number of nonterminals in the intermediate string.

Hence, the probability of seeing a finite tree is less than one!

Fortunately, most (but not all) common methods of assigning probabilities do not have this problem.

We have a "Markovian" process here (limited memory of history)

Everything above a given node in the tree is conditionally independent of everything below that node if we know what is the nonterminal in that node

Another way to think about it: once we get to a new nonterminal and continue from there, we forget the whole derivation up to that point, and focus on that nonterminal as if it is a new root node

Too strong independence assumptions for natural language data.

PCFGs and disambiguation

- A PCFG assigns a probability to every parse tree or derivation associated with a sentence.
- This probability is the product of the rules applied in building the parse tree.

$$P(T,S) = \prod_{i=1}^{n} P(A_i \to \beta_i)$$
 n is number of rules in T

- P(T,S) = P(T)P(S|T) = P(S)P(T|S) by definition
- But P(S|T) = 1 because S is determined by T

So P(T,S) = P(T)











 $P(T_{left}) = .05 * .20 * .20 * .20 * .75 * .30 * .60 * .10 * .40 = 2.2 \times 10^{-6}$ $P(T_{right}) = .05 * .10 * .20 * .15 * .75 * .30 * .60 * .10 * .40 = 6.1 \times 10^{-7}$





 $P(T_{left}) = .05 * .20 * .20 * .20 * .75 * .30 * .60 * .10 * .40 = 2.2 \times 10^{-6}$ $P(T_{right}) = .05 * .10 * .20 * .15 * .75 * .75 * .30 * .60 * .10 * .40 = 6.1 \times 10^{-7}$



 $P(T_{left}) = .05 * .20 * .20 * .20 * .75 * .30 * .60 * .10 * .40 = 2.2 \times 10^{-6}$ $P(T_{right}) = .05 * .10 * .20 * .15 * .75 * .75 * .30 * .60 * .10 * .40 = 6.1 \times 10^{-7}$



 $P(T_{left}) = .05 * .20 * .20 * .20 * .75 * .30 * .60 * .10 * .40 = 2.2 \times 10^{-6}$ $P(T_{right}) = .05 * .10 * .20 * .15 * .75 * .30 * .60 * .10 * .40 = 6.1 \times 10^{-7}$



 $P(T_{left}) = .05 * .20 * .20 * .20 * .75 * .30 * .60 * .10 * .40 = 2.2 \times 10^{-6}$ $P(T_{right}) = .05 * .10 * .20 * .15 * .75 * .30 * .60 * .10 * .40 = 6.1 \times 10^{-7}$























 $P(T_{left}) = .05 * .20 * .20 * .20 * .75 * .30 * .60 * .10 * .40 = 2.2 \times 10^{-6}$ $P(T_{right}) = .05 * .10 * .20 * .15 * .75 * .75 * .30 * .60 * .10 * .40 = 6.1 \times 10^{-7}$ Notice that shared parts have same probability!

16 / 25

Application 2: Language Modelling

As well as assigning probabilities to parse trees, a PCFG assigns a probability to every sentence generated by the grammar. This is useful for language modelling.

The probability of a sentence is the sum of the probabilities of each parse tree associated with the sentence:

$$P(S) = \sum_{Ts.t.yield(T)=S} P(T,S)$$

$$P(S) = \sum_{s.t.yield(T)=S} P(T)$$

When is it useful to know the probability of a sentence? When ranking the output of speech recognition, machine translation, and error correction systems. Many probabilistic parsers use a probabilistic version of the CYK bottom-up chart parsing algorithm.

Sentence S of length n and CFG grammar with V non-terminals

Ordinary CYK 2-d(n+1) * (n+1) array where a value in cell (i, j) is list of non-terminals spanning position *i* through *j* in *S*.

Probabilistic CYK 3-d(n+1) * (n+1) * V array where a value in cell (i, j, K) is probability of non-terminal K spanning position i through j in S

As with regular CYK, probabilistic CYK assumes that the grammar is in Chomsky-normal form (rules $A \rightarrow B \ C$ or $A \rightarrow w$).

Recursive description of probabilistic CYK

Call Chart[A, i, j] the probability of the highest-probability derivation of $w_{i+1}...w_j$ from A. Stated mathematically:

$$Chart[A, i, i + i] = p(A \to w_{i+1})$$
$$Chart[A, i, j] = \max_{\{k:i < k < j\}} \max_{\{B, C: A \to B \ C \in G\}} \max_{Chart[B, i, k] \times Chart[C, k, j] \times p(A \to B \ C)}$$

Recursive description of probabilistic CYK

Call Chart[A, i, j] the probability of the highest-probability derivation of $w_{i+1}...w_j$ from A. Stated mathematically:

$$\begin{aligned} \operatorname{Chart}[A, i, i+i] =& p(A \to w_{i+1}) \\ \operatorname{Chart}[A, i, j] = \max_{\{k:i < k < j\}} \max_{\{B, C: A \to B \ C \in G\}} \\ \operatorname{Chart}[B, i, k] \times \operatorname{Chart}[C, k, j] \times p(A \to B \ C) \end{aligned}$$

Seem familiar?

Recursive description of probabilistic CYK

Call Chart[A, i, j] the probability of the highest-probability derivation of $w_{i+1}...w_j$ from A. Stated mathematically:

$$\begin{aligned} \operatorname{Chart}[A, i, i+i] =& p(A \to w_{i+1}) \\ \operatorname{Chart}[A, i, j] = \max_{\{k:i < k < j\}} \max_{\{B, C: A \to B \ C \in G\}} \\ \operatorname{Chart}[B, i, k] \times \operatorname{Chart}[C, k, j] \times p(A \to B \ C) \end{aligned}$$

Seem familiar?

Remember CYK (where Chart[A, i, j] was simply true or false):

$$\begin{aligned} &\operatorname{Chart}[A, i, j] = &\operatorname{TRUE} \text{ iff } A \to w_{i+1} \in G \\ &\operatorname{Chart}[A, i, j] = \bigvee_{k=i+1}^{j-1} \bigvee_{A \to B} \operatorname{Chart}[B, i, k] \wedge \operatorname{Chart}[C, k, j] \end{aligned}$$

Powerful abstraction: probabilistic CYK is just CYK with a different *semiring*.

function Probabilistic-CYK(*words*, *grammar*) **returns** most probable parse and its probability

for
$$j \leftarrow$$
 from 1 to LENGTH(words) do
for all $\{A|A \rightarrow words[j] \in grammar\}$
 $table[j - 1, j, A] \leftarrow P(A \rightarrow words[j])$
for $i \leftarrow$ from $j - 2$ downto 0 do
for all $\{A|A \rightarrow BC \in grammar,$
and $table[i, k, B] > 0$ and $table[k, j, C] > 0\}$
if $(table[i, j, A] < P(A \rightarrow BC) \times table[i, k, B] \times table[k, j, C])$ then
 $table[i, j, A] \leftarrow P(A \rightarrow BC) \times table[i, k, B] \times table[k, j, C]$
 $back[i, j, A] \leftarrow \{k, B, C\}$

return

BUILD_TREE(back[1,LENGTH(words), S]), table[1,LENGTH(words), S]



	The	flight	includes	а	meal
	Det: .40				
	r1				
	[0, 1]				
		N: .02			
		[1,2]			
$S \rightarrow N$	VP VP .80) Det \rightarrow the.	40		
$NP \rightarrow$	Det N.30) $Det ightarrow a$.4	40		
$VP \rightarrow$	V NP .20	$N \rightarrow meal$	01		
$V \rightarrow i$	includes OF	$M \rightarrow flight$	12		
• /1	neruues.oc	, ingitt.			

Т	⁻ he	flight	includes	а	meal
Det	t: .40				
[0, 1	1]				
		N: .02 [1,2]			
	,		V: .05		
			[2,3]		
$S \rightarrow NP$ N	/P .80) Det \rightarrow the.	40		
NP ightarrow Det	t N.30) $Det ightarrow a$.	40		
VP ightarrow V /	NP .20	$N \to meal$.	01		
V ightarrow inclu	des.05	$N \to flight.$	02		

	The	flight	includes	а	meal
	Det: .40				
	[0, 1]				
	. /]	N: .02			
		[1,2]			
			V: .05		
			[0 3]		
			[2, 3]	Det: 40	
				DCt+0	
$S \rightarrow I$	NP VP .80	Det \rightarrow the.	40		
$NP \rightarrow$	Det N.30) Det $\rightarrow a$.	40	[3, 4]	
$VP \rightarrow$	V NP .20	$N \rightarrow meal$	01		
V ightarrow V	includes.05	$5 N \rightarrow flight.$	02		

	The	flight	includes	а	meal
	Det: .40				
	[0, 1]				
	[0, 1]	NL 00			
		N: .02			
		[1,2]			
			V: .05		
			[2, 3]		
				Det: .40	
$S \rightarrow N$	P VP .80) Det \rightarrow the.	40	I	
$NP \rightarrow$	Det N 30) Det $\rightarrow a$	40	[3, 4]	
	V ND OC	$\lambda = \Delta c c + a $			
$V r \rightarrow$	V NP .20	$N \rightarrow mean$			N: .01
$V \rightarrow ir$	ncludes.05	$N \rightarrow flight.$	02		[4,5]

	The	flight	includes	а	meal
	Det: .40	NP:			
		.30×.40 ×			
		.02 = .0024			
	[0, 1]	[0,2]			
		N: .02			
		[1,2]			
			V: .05		
			[2, 3]		
				Det: .40	
$S \rightarrow I$	VP VP .80) Det \rightarrow the.	40	[2 4]	
NP ightarrow	Det N.30) Det $ ightarrow$ a .4	40	[3, 4]	
$VP \rightarrow$	V NP .20	$N \rightarrow meal$.	01		N. 01
$V \rightarrow$	includes 05	$5 N \rightarrow flight$	02		
- / -					[4, 5]

	The	flight	includes	а	meal
	Det: .40	NP:			
		.30×.40 ×			
		.02 = .0024			
	[0, 1]	[0, 2]			
		N: .02			
		[1,2]	[1, 3]		
			V: .05		
			[2, 3]		
				Det: .40	
$S \rightarrow I$	VP VP .80) Det \rightarrow the.	40	[1 2]	
NP ightarrow	• Det N.30) Det $ ightarrow$ a .4	40	[3, 4]	
$V\!P ightarrow$	V NP .20	$N \rightarrow meal$.	01		N. 01
$V \rightarrow $	includes 05	$5 N \rightarrow flight$	12		
• / 1					[4,5]

	The	flight	includes	а	meal
	Det: .40	NP:			
		.30×.40 ×			
		.02 = .0024			
	[0, 1]	[0,2]	[0, 3]		
		N: .02			
		[1,2]	[1,3]		
			V: .05		
			[2,3]	_	
				Det: .40	
$S \rightarrow I$	VP VP .80) Det \rightarrow the.	40	[3 1]	
$NP \rightarrow$	Det N.30) Det $ ightarrow$ a .4	40	[3, 4]	
$V\!P ightarrow$	· V NP .20	N ightarrow meal .	01		N· 01
V ightarrow V	includes.05	$5 N \rightarrow flight.$	02		[4 5]
		0			[-, -]

	The	flight	includes	а	meal
	Det: .40	NP:			
		.30×.40 ×			
		.02 = .0024			
	[0, 1]	[0, 2]	[0, 3]		
		N: .02			
		[1,2]	[1,3]		
			V: .05		
			[2, 3]	[2, 4]	
				Det: .40	
$S \rightarrow I$	VP VP .80) Det \rightarrow the.	40	[2 4]	
$NP \rightarrow$	Det N.30) Det $ ightarrow$ a .4	40	[3,4]	
$V\!P ightarrow$	V NP .20	$N \rightarrow meal$.	01		N. 01
$V \rightarrow$	includes OF	$5 N \rightarrow flight 0$	02		
• / /					[4, 5]

	The	flight	includes	а	meal
	Det: .40	NP:			
		.30×.40 ×			
		.02 = .0024			
	[0, 1]	[0, 2]	[0, 3]		
		N: .02			
		[1,2]	[1, 3]	[1, 4]	
			V: .05		
			[2, 3]	[2, 4]	
				Det: .40	
$S \rightarrow I$	VP VP .80) Det \rightarrow the.	40	[2 4]	
NP ightarrow	Det N.30) Det $ ightarrow$ a .4	40	[3, 4]	
$VP \rightarrow$	V NP .20	$N \rightarrow meal$.	01		N 01
$V \rightarrow$	includes OF	5 $N \rightarrow flight$	12		IN: .U1
v -7 1	includes.00		52		[4,5]

	The	flight	includes	а	meal
	Det: .40	NP:			
		.30×.40 ×			
		.02 = .0024			
	[0, 1]	[0,2]	[0, 3]	[0, 4]	
		N: .02			
		[1,2]	[1,3]	[1, 4]	
			V: .05		
			[2, 3]	[2, 4]	
				Det: .40	
$S \rightarrow I$	VP VP .80) Det \rightarrow the.	40	[3 4]	
$NP \rightarrow$	• Det N.30) Det $ ightarrow$ a .4	40	[3, 4]	
$V\!P ightarrow$	· V NP .20	N ightarrow meal .	01		N· 01
V ightarrow V	includes.05	$5 N \rightarrow flight.$	02		[4 5]
		0			[-[-, -]

	The	flight	includes	а	meal
	Det: .40	NP:			
		.30×.40 ×			
		.02 = .0024			
	[0, 1]	[0, 2]	[0, 3]	[0, 4]	
		N: .02			
		[1,2]	[1,3]	[1, 4]	
			V: .05		
			[2, 3]	[2, 4]	
			Det: .40		
					$NP: .30 \times .40 \times .01 = 0.0010$
$S \rightarrow NP VP$.80 $Det \rightarrow the$.40					.01 = 0.0012
$NP \rightarrow Det N.30$ $Det \rightarrow a$.40				[3, 4]	[3, 5]
$VP \rightarrow V NP 20 N \rightarrow meal 01$					
$V \rightarrow includes 05 N \rightarrow flight 02$					N: .01
$v \rightarrow \text{Includes.05}$ $N \rightarrow \text{Tight.02}$					[4,5]

	The flight		includes	а	meal
	Det: .40	NP:			
		.30×.40 ×			
		.02 = .0024			
	[0,1]	[0,2]	[0, 3]	[0, 4]	
		N: .02			
		[1,2]	[1, 3]	[1, 4]	
		V: .05		VP: .20 ×	
				$.05 \times 0.0012 =$	
				0.000012	
		[2, 3]	[2, 4]	[2,5]	
			Det: .40	ND 20 40	
					NP: $.30 \times .40 \times$
$S \rightarrow NP \ VP \ .80 Det \rightarrow the.40$				[0, 4]	.01 = 0.0012
$NP \rightarrow Det N.30$ $Det \rightarrow a$.40				[3, 4]	[3, 5]
$VP \rightarrow V NP .20 N \rightarrow meal .01$					N 01
$V \rightarrow includes 05 N \rightarrow flight 02$					N: .01
$v \rightarrow \text{Includes.05}$ $N \rightarrow \text{Inglit.02}$					[4, 5]

	The flight		includes	а	meal
	Det: .40	NP:			
		.30×.40 ×			
		.02 = .0024			
	[0, 1]	[0,2]	[0, 3]	[0, 4]	
		N: .02			
[1,2]		[1,2]	[1, 3]	[1, 4]	[1,5]
		V: .05		VP: .20 ×	
				$.05 \ imes \ 0.0012 \ =$	
				0.000012	
		[2, 3]	[2, 4]	[2,5]	
			Det: .40	ND 00 40	
					NP: $.30 \times .40 \times$
$S \rightarrow NP VP$.80 $Det \rightarrow the.40$				[o 4]	.01 = 0.0012
$NP \rightarrow Det N.30$ $Det \rightarrow a$.40				[3, 4]	[3, 5]
$VP \rightarrow V NP 20 N \rightarrow meal 01$					
$V \rightarrow includes 05 N \rightarrow flight 02$					N: .01
$V \rightarrow \text{Includes.05}$ $N \rightarrow \text{flight.02}$					[4,5]

	The flight		includes	а	meal
	Det: .40 NP: .30×.40 ×				S: .80 \times .0024 \times
					.000012 =
		.02 = .0024			.00000023
	[0, 1]	[0, 2]	[0, 3]	[0, 4]	[0, 5]
		N: .02			
[1,2]		[1,2]	[1,3]	[1, 4]	[1,5]
		V: .05		VP: .20 ×	
				$.05 \ imes \ 0.0012 \ =$	
				0.000012	
		[2, 3]	[2, 4]	[2,5]	
			Det: .40	ND 20 40	
					NP: $.30 \times .40 \times$
$S \rightarrow NP VP .80 Det \rightarrow the.40$.01 = 0.0012
$NP \rightarrow Det N.30$ $Det \rightarrow a$.40				[3, 4]	[3, 5]
$VP \rightarrow V NP 20 N \rightarrow meal 01$					
V_{i} is alreaded OF N_{i} if the the OP					N: .01
$V \rightarrow includes.05 N \rightarrow flight.02$				[4,5]	

Probabilistic CYK: more tricky example



(Not quite in CNF, but never mind.) We'll parse:

orange tree blossoms early

The probabilistic CYK-style chart

	orange	tree	blossoms	early
orange	N (0.3)	NP (0.06)	S (0.048)	S (0.012)
	A (1.0)		NP (0.0024)	
	NP (0.18)			
tree		N (0.5)	NP (0.012)	S (0.06)
		NP (0.3)		
blossoms			N (0.2)	VP (0.2)
			V (1.0)	
			NP (0.12)	
			VP (0.8)	
early				Adv(1.0)

• The phrase orange tree gets 0.06 for its best analysis *as an NP*, since

 $\begin{array}{rl} 0.06 &= 0.2*1.0*0.3 & (\mbox{for NP} \rightarrow A \mbox{ NP}) \\ \mbox{beats } 0.018 &= 0.18*0.5*0.2 & (\mbox{for NP} \rightarrow NP \mbox{ N}). \\ \mbox{Only the higher probability is recorded in the chart.} \end{array}$

- For orange tree blossoms, there are now two analyses as NP, each with probability 0.0024.
- There is also an analysis of orange tree blossoms as S. This doesn't compete with its analysis as NP, so both are recorded.

- A PCFG is a CFG with each rule annotated with a probability;
- the sum of the probabilities of all rules that expand the same non-terminal must be 1;
- probability of a parse tree is the product of the probabilities of all the rules used in this parse;
- probability of sentence is sum of probabilities of all its parses;
- applications for PCFGs: disambiguation, language modeling;
- Probabilistic CYK algorithm.

Next lecture: But where do the rule probabilities come from?