

1 ICL/PCFGs/2005–11–17

Contents

2 Probabilistic context-free grammars

Context-free Grammars (CFGs)

Recall...

A CFG is a 4-tuple $\langle N, \Sigma, P, S \rangle$ where:

- N is the set of non-terminal symbols (eg S, NP, VP, \dots)
- Σ is the set of terminal symbols (eg *the, cat, sat, \dots*)
- P is the set of productions $A \rightarrow \alpha$, $A \in N$ and $\alpha \in (N \cup \Sigma)^*$ (eg $S \rightarrow VP NP$)
- S is the start symbol S

Parsing

- *Top-down* Predict what you expect to see (eg Earley algorithm)
- *Bottom-up* Start with the words, then incrementally build up parse trees
 - CYK (Cocke-Younger-Kasami) algorithm
 - Well matched to probabilistic grammars—assign probabilities to constituents as they are completed and placed in the table
 - Dynamic programming approach—store and reuse intermediate results (ie probabilities of sub-trees)
 - Resolve ambiguities by taking the most probable sub-tree

Chomsky Normal Form (CNF)

- If a CFG is in *Chomsky Normal Form*, productions are constrained to be of two forms only:
 - **Expand to 2 non-terminals**, eg: $A \rightarrow BC$
with A, B, C all non-terminals
 - **Expand to 1 terminal**, eg: $A \rightarrow a$
where A is a non-terminal and a is a terminal
- Any CFG can be translated to a (weakly) equivalent CFG in CNF
- The CYK algorithm requires a grammar in CNF

3 Bottom-up parsing

CYK Example

Alice called Bob from Cardiff

Context-Free Grammar:

$S \rightarrow NP VP$ $NP \rightarrow Alice$
 $NP \rightarrow NP PP$ $NP \rightarrow Bob$
 $VP \rightarrow V NP$ $NP \rightarrow Cardiff$
 $VP \rightarrow VP PP$ $V \rightarrow called$
 $PP \rightarrow P NP$ $P \rightarrow from$

CYK Parsing

				NP
			P	Cardiff
		NP	from	
	V	Bob		
NP	called			
Alice				

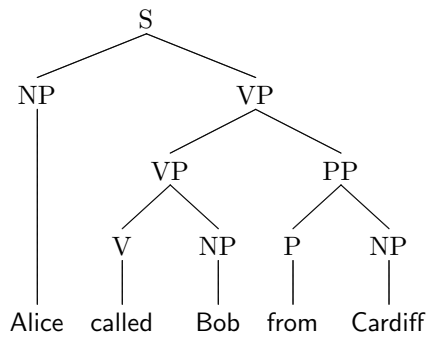
CYK Parse Chart

CYK Parsing

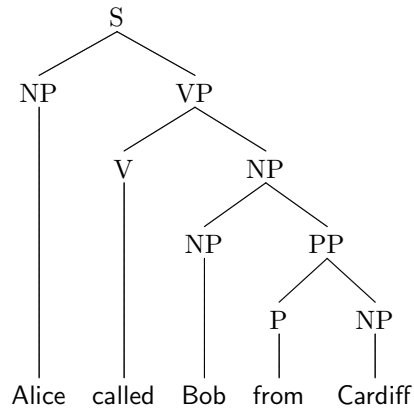
S	VP₁	NP	PP	NP
X	X	X	P	Cardiff
S	VP	NP	from	
X	V	Bob		
NP	called			
Alice				

First Parse

First tree



Second tree



4 Probabilistic CYK

Probabilistic context-free grammars (PCFGs)

- A probabilistic context-free grammar augments each rule in a CFG with a conditional probability p
 $A \rightarrow \alpha \quad (p)$
- This probability is the probability that given non-terminal A it will be expanded to the sequence α ;
 written as $P(A \rightarrow \alpha|A)$ or $P(A \rightarrow \alpha)$
- Probability of a parse tree is the product of the rule probabilities used to construct the parse

Probabilistic parsing

- Consider the rule:

$$S \rightarrow NP VP$$

Then the probability of S is the product of the rule probability and the probability of each subtree:

$$P(S) = P(S \rightarrow NP VP) \cdot P(NP) \cdot P(VP)$$

- We are doing bottom-up parsing... so we already know the subtree probabilities $P(NP)$ and $P(VP)$

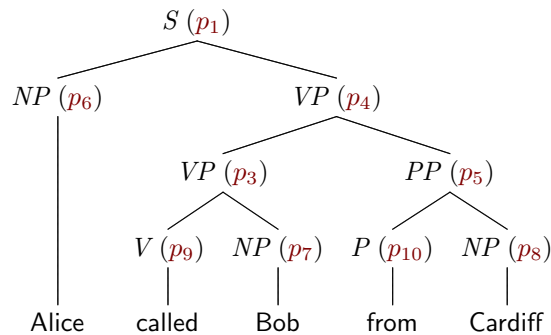
Probabilistic CYK Example

Alice called Bob from Cardiff

Context-Free Grammar:

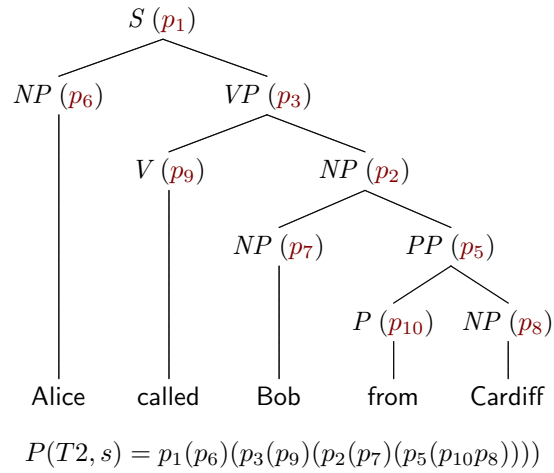
$S \rightarrow NP VP$	(p_1)	$NP \rightarrow \text{Alice}$	(p_6)
$NP \rightarrow NP PP$	(p_2)	$NP \rightarrow \text{Bob}$	(p_7)
$VP \rightarrow V NP$	(p_3)	$NP \rightarrow \text{Cardiff}$	(p_8)
$VP \rightarrow VP PP$	(p_4)	$V \rightarrow \text{called}$	(p_9)
$PP \rightarrow P NP$	(p_5)	$P \rightarrow \text{from}$	(p_{10})

First tree (T1)



$$P(T1, S) = p_1(p_6)(p_4(p_3(p_9p_7))(p_5(p_{10}p_8)))$$

Second tree (T2)



Choosing the tree

Choose tree 1 if $P(T1, S)/P(T2, S) > 1$.

$$\frac{P(T1, S)}{P(T2, S)} = \frac{p_1(p_6)(p_4(p_3(p_9p_7))(p_5(p_{10}p_9)))}{p_1(p_6)(p_3(p_9)(p_2(p_7)(p_5(p_{10}p_8))))} = \frac{p_4}{p_2}$$

Probabilistic CYK

If $p_4p_3p_7p_9p_5p_8p_{10} > p_3p_9p_2p_7p_5p_8p_{10}$:

<i>S</i> $p_1p_6p_4p_3p_7p_9p_5p_8p_{10}$	<i>VP₂</i> $p_3p_9p_2p_7p_5p_8p_{10}$	<i>NP</i> $p_2p_7p_5p_8p_{10}$	<i>PP</i> $p_5p_8p_{10}$	<i>NP</i> p_8
X	X	X	<i>P</i> p_{10}	Cardiff
S $p_1p_6p_3p_7p_9$	<i>VP</i> $p_3p_7p_9$	<i>NP</i> p_7	from	
X	<i>V</i> p_9	Bob		
<i>NP</i> p_6	called			
Alice				

5 Training PCFGs

Estimating PCFG Probabilities

- Treebank—corpus of parsed sentences
- Given a treebank compute the probability of each non-terminal expansion ($A \rightarrow \alpha$) based on the counts $c(A \rightarrow \alpha)$:

$$P(A \rightarrow \alpha|A) = \frac{c(A \rightarrow \alpha)}{\sum_Y c(A \rightarrow Y)} = \frac{c(A \rightarrow \alpha)}{c(A)}$$

- If a treebank is not available probabilities estimated by parsing the corpus, incrementing a counter for each rule in the parse, and normalizing to get probabilities

- Problem: multiple possible parses for most sentences
- Solution: keep separate counts for each parse, and weight by the probability of the parse. Can be computed using the Inside-Outside algorithm.

6 Probabilistic lexicalised CFGs

Problems with PCFGs

Structural Rule probabilities are independent of location in the parse tree. For example pronouns are much more likely to be subjects than objects

Lexical PCFGs only capture lexical information in the expansion of pre-terminals.

But, lexical dependencies can often be used to choose the correct parse, eg:

Carol eats chips with ketchup

Carol eats chips with a fork

Lexical dependence

Data from Penn Treebank:

Rule	<i>come</i>	<i>take</i>	<i>think</i>	<i>want</i>
VP → V	0.095	0.026	0.046	0.057
VP → V NP	0.011	0.321	0.002	0.139
VP → V PP	0.345	0.031	0.071	0.003
VP → V S	0.022	0.013	0.048	0.708

The rule used to expand VP is strongly dependent on the verb

Lexicalised PCFGs

- Annotated each non-terminal with its *lexical head*
- Each rule has a head child on the left hand side; the headword for a node is then the headword of its head child
- Easy for simple examples (eg the noun is the head of an NP, the verb is the head of a VP); harder in practice
- Various heuristics for finding the head robustly (mainly developed on Penn Treebank)
- A “simple” lexicalised CFG is a basic CFG with a lot more rules (ie each rule is copied for each headword) — but this is impractical!
- Make some independence assumptions and condition the probability of each rule on the head (or nearby heads)

7 Summary

- **Reading:** J+M, chapter 12 (12.1 - 12.3)
- Bottom-up CYK parsing of CFGs
- Probabilistic CFGs and probabilistic parsing using CYK
- Lexical and structural problems with PCFGs; lexicalised PCFGs.