

Probabilistic context-free parsing

Parsing
ISCL-BA-06

Çağrı Çöltekin
ccolt@informatik.uni-tuebingen.de

University of Tübingen
Seminar für Sprachwissenschaft

Winter Semester 2020/21

www.informatik.uni-tuebingen.de

Context-free grammars

recap

- Context free (CF) grammars are most practically useful grammars in the Chomsky hierarchy
- Most of the parsing theory (and practice) is build on parsing CF languages
- The context-free rules have the form

$$A \rightarrow \alpha$$

where A is a single non-terminal symbol and α is a (possibly empty) sequence of terminal or non-terminal symbols

C. Çöltekin, INF | University of Tübingen

Winter Semester 2020/21 1 / 30

An example context-free grammar

Derivation of sentence 'she saw a duck'

$S \rightarrow NP VP$
 $S \rightarrow \text{Aux NP VP}$
 $NP \rightarrow \text{Det N}$
 $NP \rightarrow \text{Prn}$
 $NP \rightarrow NP PP$
 $VP \rightarrow V NP$
 $VP \rightarrow V$
 $VP \rightarrow VP PP$
 $PP \rightarrow \text{Prep NP}$
 $N \rightarrow \text{duck}$
 $N \rightarrow \text{park}$
 $N \rightarrow \text{parke}$
 $V \rightarrow \text{duck}$
 $V \rightarrow \text{dacks}$
 $V \rightarrow \text{saw}$
 $Prn \rightarrow \text{she / her}$
 $Prep \rightarrow \text{in / with}$
 $Det \rightarrow \text{a / the}$

C. Çöltekin, INF | University of Tübingen

Winter Semester 2020/21 2 / 30

Representations of a context-free parse tree

A parse tree

A history of derivations:

- $S \rightarrow NP VP$
- $NP \rightarrow Prn$
- $Prn \rightarrow I$
- $VP \rightarrow V NP$
- $V \rightarrow \text{saw}$
- $NP \rightarrow Det N$
- $Prn_{12} \rightarrow \text{her}$
- $N \rightarrow \text{duck}$

A sequence with (label) brackets

$$\left[\left[\left[\text{NP } \text{Prn } I \right] \right] \left[\text{VP } \text{V } \text{saw} \left[\left[\text{NP } \text{Det } \text{a } \text{N } \text{duck} \right] \right] \right]$$

C. Çöltekin, INF | University of Tübingen

Winter Semester 2020/21 3 / 30

Parsing with context-free grammars

- Parsing can be
 - top-down: start from S , search for derivation that leads to the input
 - bottom-up: start from input, try to reduce it to S
- Naïve search for both recognition/parsing is intractable
- Dynamic programming methods allow polynomial time recognition
 - CKY bottom-up, requires Chomsky normal form
- Early top-down (with bottom-up filtering), works with unrestricted grammars
 - $O(n^3)$ time complexity (for recognition)
- Chart parsers are (reasonably) efficient, and they can represent ambiguity in their output
- However, they do not help with resolving ambiguity

C. Çöltekin, INF | University of Tübingen

Winter Semester 2020/21 4 / 30

Natural languages are ambiguous



C. Çöltekin, INF | University of Tübingen

Winter Semester 2020/21 5 / 30

Some types of ambiguities

- Lexical ambiguity
 - She is looking for a match
 - We saw her duck
- Attachment ambiguity
 - I saw the man with a telescope
 - Panda eats bamboo shoots and leaves
- Local ambiguity (garden path sentences)
 - The horse raced past the barn fell
 - The old man the boats
 - Fat people eat accumulates

C. Çöltekin, INF | University of Tübingen

Winter Semester 2020/21 6 / 30

Ambiguity and the parsers

- Given a grammar, chart parsers (e.g., CKY, Early) can parse natural language sentences relatively efficiently
- These parsers also represent all possible parse trees in their chart/output efficiently
- However, they have nothing to say about which of these parses are the most likely one.
- The task of selecting the best parse among many is called disambiguation
- In almost all practical uses, parsers are combined with disambiguators

C. Çöltekin, INF | University of Tübingen

Winter Semester 2020/21 7 / 30

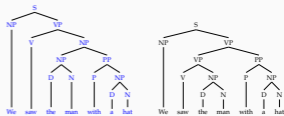
We do not recognize many ambiguities

- Time flies like an arrow; fruit flies like a banana
 - Outside of a dog, a book is a man's best friend; inside it's too hard to read
 - One morning I shot an elephant in my pajamas. How he got in my pajamas, I don't know.
 - Don't eat the pizza with a knife and fork; the one with mushrooms is better.
- A parser, nevertheless, produces multiple parses for these sentences.

C. Çöltekin, INF | University of Tübingen

Winter Semester 2020/21 8 / 30

The task: choosing the most plausible parse



C. Çöltekin, INF | University of Tübingen

Winter Semester 2020/21 9 / 30

Statistical parsing

- Find the most plausible parse of an input string given all possible parses
- We need a scoring function, for each parse, given the input
- We typically use probabilities for scoring, task becomes finding the parse (or tree), t , given the input string w

$$t_{\text{best}} = \underset{t}{\text{arg max}} P(t | w)$$

- Note that some ambiguities need a larger context than the sentence to be resolved correctly

C. Çöltekin, INF | University of Tübingen

Winter Semester 2020/21 10 / 30

Probability refresher (1)

- Probability is a measure of (un)certainly of an event
- We quantify the probability of an event with a number between 0 and 1
 - 0 the event is impossible
 - 0.5 the event is as likely to happen (or happened) as it is not
 - 1 the event is certain
- All possible outcomes of a trial (experiment or observation) is called the sample space (Ω)

Axioms of probability states that

- $P(E) \in \mathbb{R}, P(E) \geq 0$
- $P(\Omega) = 1$
- For disjoint events E_1 and E_2 , $P(E_1 \cup E_2) = P(E_1) + P(E_2)$

C. Çöltekin, INF | University of Tübingen

Winter Semester 2020/21 11 / 30

Probability refresher (2)

Joint and conditional probabilities, chain rule

- Joint probability of two events is noted as $P(x, y)$

- The conditional probability is defined as

$$P(x|y) = \frac{P(x, y)}{P(y)} \text{ or } P(x, y) = P(x|y)P(y)$$

- If the events x and y are independent,

$$P(x|y) = P(x), P(y|x) = P(y), P(x, y) = P(x)P(y)$$

- For more than two variables (chain rule):

$$P(x, y, z) = P(z, y)P(y|x) = P(x|y, z)P(y|z)P(z) = \dots$$

- If all are independent

$$P(x, y, z) = P(x)P(y)P(z)$$

Probabilistic context free grammars (PCFG)

- A probabilistic context free grammar augments a CFG by adding a probability value to each rule

$$A \rightarrow \alpha \quad |p|$$

where A is a non-terminal, α is string of terminals and non-terminals, and p is the probability associated with the rule

- Like CFGs, a PCFG accepts a sentence if it can be derived from S with rules $R_1 \dots R_k$

- The probability of a parse tree t of input string w , $P(t|w)$, corresponding to the derivation $R_1 \dots R_k$ is

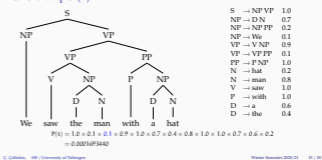
$$P(t|w) = \prod_{i=1}^k p(R_i)$$

where $p(R_i)$ is the probability of the rule R_i

PCFG example (1)



PCFG example (2)



Where do the rule probabilities come from?

- Supervised: estimate from a treebank, e.g., using maximum likelihood estimation
- Unsupervised: expectation-maximization (EM)

PCFGs - an interim summary

- PCFGs assign probabilities to parses based on CFG rules used during the parse
- PCFGs assume that the rules are independent
- PCFGs are generative models, they assign probabilities to $P(t, w)$, we can calculate the probability of a sentence by

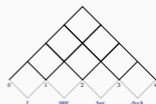
$$P(w) = \sum_t P(t, w) = \sum_t P(t)$$

PCFG chart parsing

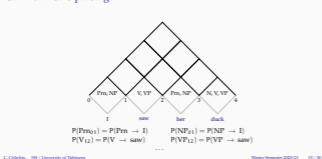
- Both CKY and Earley algorithms can be adapted to PCFG parsing
- CKY matches PCFG parsing quite well
 - to get the best parse, store the constituent with the highest probability in every cell of the chart
 - to get n -best parse (beam search), store the n -best constituents in every cell in the chart

CKY for PCFG parsing

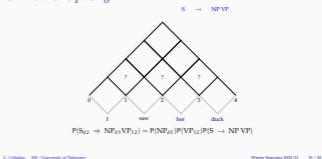
CFG essay: Ambiguity, Statistical parsing, PCFGs, Evaluation



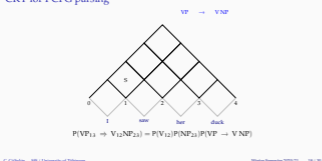
CKY for PCFG parsing



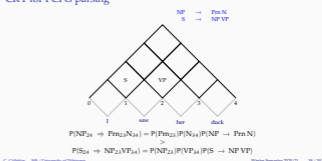
CKY for PCFG parsing



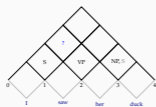
CKY for PCFG parsing



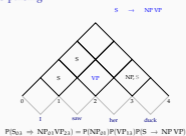
CKY for PCFG parsing



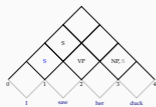
CKY for PCFG parsing



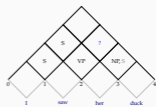
CKY for PCFG parsing



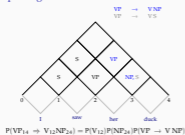
CKY for PCFG parsing



CKY for PCFG parsing



CKY for PCFG parsing



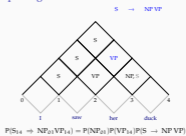
CKY for PCFG parsing



CKY for PCFG parsing



CKY for PCFG parsing



CKY for PCFG parsing



CKY for PCFG parsing



What makes the difference in PCFG probabilities?

S \rightarrow NP VP	1.0	S \rightarrow NP VP	1.0
NP \rightarrow We	0.1	NP \rightarrow We	0.1
VP \rightarrow VP PP	0.1	VP \rightarrow V NP	0.7
VP \rightarrow V NP	0.8	V \rightarrow saw	1.0
V \rightarrow saw	1.0	NP \rightarrow NP PP	0.2
NP \rightarrow D N	0.7	NP \rightarrow D N	0.7
D \rightarrow the	0.4	D \rightarrow the	0.4
N \rightarrow man	0.8	N \rightarrow man	0.8
PP \rightarrow P NP	1.0	PP \rightarrow P NP	1.0
P \rightarrow with	1.0	P \rightarrow with	1.0
NP \rightarrow D N	0.7	NP \rightarrow D N	0.7
D \rightarrow a	0.6	D \rightarrow a	0.6
N \rightarrow hat	0.2	N \rightarrow hat	0.2

The parser's choice would not be affected by lexical items!

What is wrong with PCFGs?

- In general: the assumption of independence
- The parents affect the correct choice for children, for example, in English NP \rightarrow I'm is more likely in the subject position
- The lexical units affect the correct decision, for example:
 - We eat the pizza with hands
 - We eat the pizza with mushrooms
- Additionally: PCFGs use local context, difficult to incorporate arbitrary/global features for disambiguation

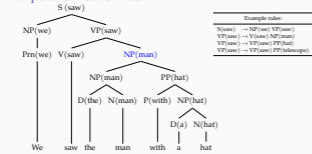
Solutions to PCFG problems

- Independence assumptions can be relaxed by either
 - Parent annotation
 - Lexicalization
 - Re-ranking
- To condition on arbitrary/global information: discriminative models
- Most practical PCFG parsers are lexicalized, and often use a re-ranker conditioning on other (global) features

Lexicalizing PCFGs

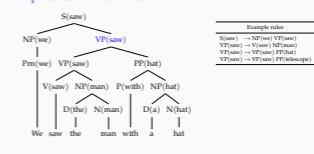
- Replace non-terminal X with $X(h)$, where h is a tuple with the lexical word and its POS tag
- Now the grammar can capture (head-driven) lexical dependencies
- But number of nonterminals grow by $|V| \times |T|$
- Estimation becomes difficult (many rules, data sparsity)
- Some treebanks (e.g., Penn Treebank) do not annotate heads, they are automatically annotated (based on heuristics)

Example lexicalized derivation



Example rules:	
$S(\text{saw})$	$\rightarrow NP(\text{we}) VP(\text{saw})$
$VP(\text{saw})$	$\rightarrow V(\text{saw}) NP(\text{man})$
$NP(\text{man})$	$\rightarrow D(\text{the}) N(\text{man})$
$NP(\text{man})$	$\rightarrow VP(\text{with}) PP(\text{hat})$
$VP(\text{with})$	$\rightarrow V(\text{with}) PP(\text{hat})$

Example lexicalized derivation



Example rules:	
$S(\text{saw})$	$\rightarrow NP(\text{we}) VP(\text{saw})$
$VP(\text{saw})$	$\rightarrow V(\text{saw}) NP(\text{man})$
$VP(\text{saw})$	$\rightarrow VP(\text{with}) PP(\text{hat})$
$VP(\text{with})$	$\rightarrow V(\text{with}) PP(\text{hat})$

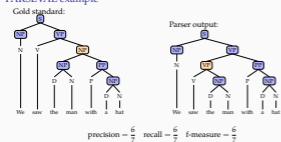
Evaluating the parser output

- A parser can be evaluated extrinsically based on its effect on a task (e.g., machine translation) where it is used
- Intrinsically based on the match with ideal parsing
- The typically evaluation (intrinsic) is based on a gold standard (GS)
- Exact match is often
 - very difficult to achieve (think about a 50-word newspaper sentence)
 - not strictly necessary (recovering parts of the parse can be useful for many purposes)

Parser evaluation metrics

- Common evaluation metrics are (PARSEVAL):
 - precision: the ratio of correctly predicted nodes
 - recall: the nodes (in GS) that are predicted correctly
 - f -measure: harmonic mean of precision and recall $\left(\frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}} \right)$
- The measures can be
 - unlabeled: the spans of the nodes are expected to match
 - labelled: the node label should also match
- Crossing brackets (or average non-crossing brackets)
 - (We (saw (them (with (binoculars)))
 - (We ((saw them) (with binoculars)))
- Measures can be averaged per constituent (micro average), or over sentences (macro average)

PARSEVAL example



Problems with PARSEVAL metrics

- PARSEVAL metrics favor certain type of structures
 - Results are surprisingly well for flat tree structures (e.g., Penn treebank)
 - Results of some mistakes are catastrophic (e.g., low attachment)
- Not all mistakes are equally important for semantic distinctions
- Some alternatives:
 - Extrinsic evaluation
 - Evaluation based on extracted dependencies

Summary

- PCFG is a simple attempt to augment CFG with probabilities
 - PCFG parsing alone is suboptimal: Independence assumptions are too strong
 - Solutions include (a combination of) lexicalization, parent annotation and re-ranking
 - Reading suggestion: Jurafsky and Martin (2009, Chapter 14)
- Next:
- Dependency grammars and dependency parsing

Acknowledgments, references, additional reading material

- Steve Chalk and David J. Swales (2007). *Forming Dependency & Practical Guide, second ed. Morphology in Computer Science: Theoretical and Applied Aspects (Morphology and Natural Language Processing and Dependency Grammar)*. Springer, New York, ISBN 978-0-387-08449-0
- Andrew Senior and James R. Martin (2019). *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition, second ed. Second Edition (ISBN: 978-0-07-076503-0)*. <http://nlp.stanford.edu/~jrs/pubs/1919/>

