

- 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

$$A \rightarrow \alpha$$

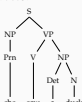
where A is a single non-terminal symbol and α is a (possibly empty) sequence of terminal or non-terminal symbols

An example context-free grammar

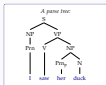
S \rightarrow NP VP
S \rightarrow Aux NP VP
NP \rightarrow Det N
NP \rightarrow Prn
NP \rightarrow NP PP
VP \rightarrow V NP
VP \rightarrow V
VP \rightarrow VP PP
PP \rightarrow Prep NP
N \rightarrow duck
N \rightarrow park
N \rightarrow parks
V \rightarrow duck
V \rightarrow ducks
V \rightarrow saw
Prn \rightarrow she / her
Prep \rightarrow in (with)
Det \rightarrow a / the

Derivation of sentence 'she saw a duck'

S \rightarrow NP VP
NP \rightarrow Prn
Prn \rightarrow she
VP \rightarrow V NP
V \rightarrow saw
NP \rightarrow Det N
Det \rightarrow a
N \rightarrow duck



Representations of a context-free parse tree



A history of derivations:

- S \rightarrow NP VP
- NP \rightarrow Prn
- Prn \rightarrow I
- VP \rightarrow V NP
- V \rightarrow saw
- NP \rightarrow Prn, N
- Prn \rightarrow her
- N \rightarrow duck

A sequence with (labeled) brackets

$\left[\left[\left[\text{I} \right]_{\text{NP}} \left[\text{saw} \right]_{\text{VP}} \right]_{\text{VP}} \left[\left[\left[\text{her} \right]_{\text{NP}} \left[\text{duck} \right]_{\text{N}} \right]_{\text{NP}} \right]_{\text{NP}} \right]$

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
- Naive 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

Natural languages are ambiguous



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

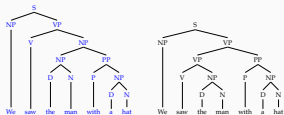
Ambiguity and the parsers

- 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

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.

The task: choosing the most plausible parse



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

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

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

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 , $P(E_1 \cup E_2) = P(E_1) + P(E_2)$

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):

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

- 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)

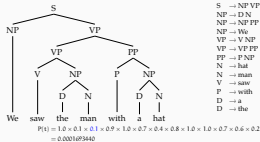
CFG essay: Ambiguity Statistical parsing PCFGs Evaluation



$S \rightarrow NP VP$	1.0
$NP \rightarrow D N$	0.7
$NP \rightarrow NP PP$	0.2
$NP \rightarrow We$	0.1
$VP \rightarrow V NP$	0.9
$VP \rightarrow VP PP$	0.1
$PP \rightarrow P NP$	1.0
$N \rightarrow hat$	0.2
$N \rightarrow man$	0.8
$N \rightarrow man$	0.8
$V \rightarrow saw$	1.0
$P \rightarrow with$	1.0
$D \rightarrow a$	0.6
$D \rightarrow the$	0.4

PCFG example (2)

CFG essay: Ambiguity Statistical parsing PCFGs Evaluation



$S \rightarrow NP VP$	1.0
$NP \rightarrow D N$	0.7
$NP \rightarrow NP PP$	0.2
$NP \rightarrow We$	0.1
$VP \rightarrow V NP$	0.9
$VP \rightarrow VP PP$	0.1
$PP \rightarrow P NP$	1.0
$N \rightarrow hat$	0.2
$N \rightarrow man$	0.8
$N \rightarrow man$	0.8
$V \rightarrow saw$	1.0
$P \rightarrow with$	1.0
$D \rightarrow a$	0.6
$D \rightarrow the$	0.4

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

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

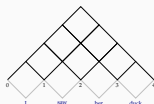
PCFG chart parsing

CFG essay: Ambiguity Statistical parsing PCFGs Evaluation

- 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 best parse (beam search), store the n -best constituents in every cell in the chart

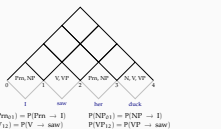
CKY for PCFG parsing

CFG essay: Ambiguity Statistical parsing PCFGs Evaluation



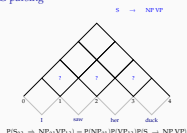
CKY for PCFG parsing

CFG essay: Ambiguity Statistical parsing PCFGs Evaluation



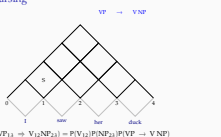
CKY for PCFG parsing

CFG essay: Ambiguity Statistical parsing PCFGs Evaluation



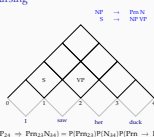
CKY for PCFG parsing

CFG essay: Ambiguity Statistical parsing PCFGs Evaluation

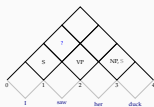


CKY for PCFG parsing

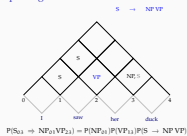
CFG essay: Ambiguity Statistical parsing PCFGs Evaluation



CKY for PCFG parsing

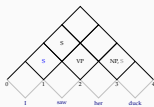


CKY for PCFG parsing

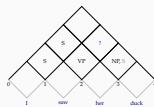


$$P(S_{01} \Rightarrow NP_0(VP_{21}) = P(NP_{01})P(VP_{11})P(S \rightarrow NP VP)$$

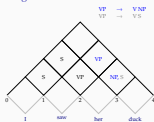
CKY for PCFG parsing



CKY for PCFG parsing

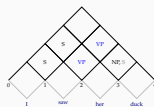


CKY for PCFG parsing

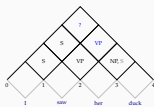


$$P(VP_{14} \Rightarrow V_{12}NP_{24}) = P(V_{12})P(NP_{24})P(VP \rightarrow V NP)$$

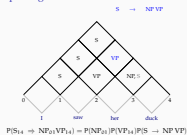
CKY for PCFG parsing



CKY for PCFG parsing

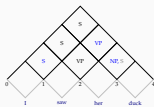


CKY for PCFG parsing

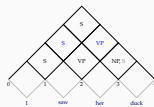


$$P(S_{14} \Rightarrow NP_0(VP_{14}) = P(NP_{01})P(VP_{14})P(S \rightarrow NP VP)$$

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 P'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

