

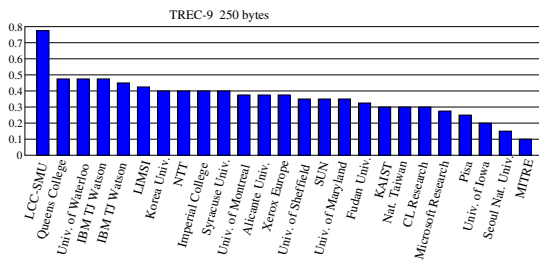
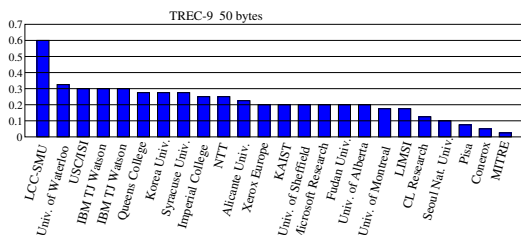
{Probabilistic|Stochastic}

Context-Free Grammars (PCFGs)

FSNLP, chapter 11

Christopher Manning and
Hinrich Schütze
© 1999–2002

301



303

Question Answering Example (1)

- Q261: What company sells most greetings cards ?
- sells(ORGANIZATION, cards(greeting), most)
- "Hallmark remains the largest maker of greeting cards" maker(ORGANIZATION(Hallmark), cards(greeting), largest)
- Need an entailment between *producing*, or *making* and *selling goods*
- Derived from WordNet, since synset *make*, *produce*, *create* has the genus **manufacture**, defined in the gloss of its nominalization as (*for*) *sale*
- Also, need *most* ≈ *largest*
- Therefore the semantic form of question Q261 and its illustrated answer are similar

305

Question answering from text

- TREC 8+ QA competition (1999–; it's ongoing): an idea originating from the IR community
- With massive collections of on-line documents, manual translation of textual information into knowledge bases covering large numbers of domains is impractical: We want to answer questions from textbases
- Evaluated output is 5 answers of 50/250 byte snippets of text drawn from a 3 Gb text collection, and required to contain at least one concept of the semantic category of the expected answer type. (Until 2002. IR think: Suggests the use of named entity recognizers.)
- Get reciprocal points for highest correct answer.

302

Pasca and Harabagiu (2001) demonstrates the value of sophisticated NLP processing

- Good IR is needed: paragraph retrieval based on SMART
- Large taxonomy of question types and expected answer types is crucial
- Parsing: A statistical parser (modeled on Collins 1997) is used to parse questions, relevant text for answers, and WordNet to build a knowledge base for reasoning
- Controlled query expansion loops (morphological, lexical synonyms, and semantic relations) are all important in retrieving the correct answer.
- Answer ranking by ML method based on this information surpasses IR-style empirical methods.

304

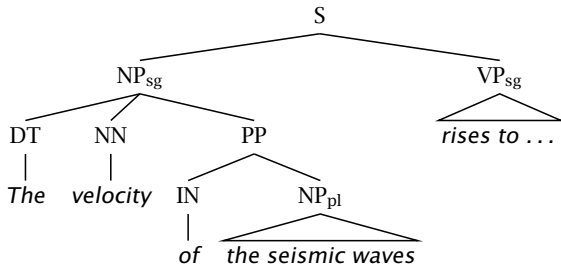
Question Answering Example (2)

- How hot does the inside of an active volcano get ?
- get(TEMPERATURE, inside(volcano(active)))
- "lava fragments belched out of the mountain were as hot as 300 degrees Fahrenheit"
- fragments(lava, TEMPERATURE(degrees(300)), belched(out, mountain))
- □ volcano ISA mountain
- □ lava ISPARTOF volcano □ lava inside volcano
- □ fragments of lava HAVEPROPERTIESOF lava
- The needed semantic information is available in WordNet definitions, and was successfully translated into a form that can be used for rough 'proofs'

306

Why we need recursive phrase structure

- *The velocity of the seismic waves rises to ...*
- Kupiec (1992): HMM tagger goes awry: waves → verb



307

PCFGs

A PCFG G consists of the usual parts of a CFG

- A set of terminals, $\{w^k\}, k = 1, \dots, V$
- A set of nonterminals, $\{N^i\}, i = 1, \dots, n$
- A designated start symbol, N^1
- A set of rules, $\{N^i \rightarrow \zeta^j\}$, (where ζ^j is a sequence of terminals and nonterminals)

and

- A corresponding set of probabilities on rules such that:

$$\forall i \sum_j P(N^i \rightarrow \zeta^j) = 1$$

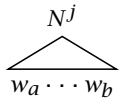
308

PCFG notation

Sentence: sequence of words $w_1 \dots w_m$

w_{ab} : the subsequence $w_a \dots w_b$

N_{ab}^i : nonterminal N^i dominates $w_a \dots w_b$



$N^i \xRightarrow{*} \zeta$: Repeated derivation from N^i gives ζ .

309

PCFG probability of a string

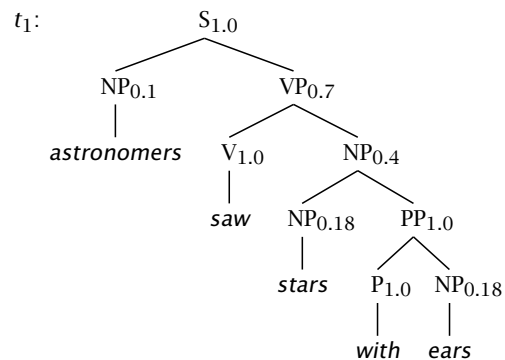
$$\begin{aligned} P(w_{1n}) &= \sum_t P(w_{1n}, t) \quad t \text{ a parse of } w_{1n} \\ &= \sum_{\{t: \text{yield}(t)=w_{1n}\}} P(t) \end{aligned}$$

310

A simple PCFG (in CNF)

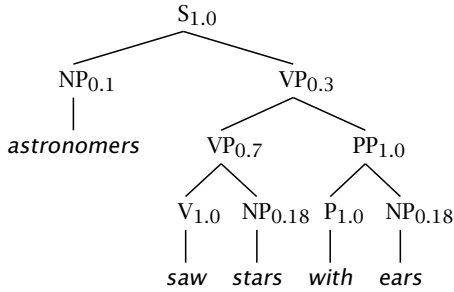
$S \rightarrow NP VP$	1.0	$NP \rightarrow NP PP$	0.4
$PP \rightarrow P NP$	1.0	$NP \rightarrow \textit{astronomers}$	0.1
$VP \rightarrow V NP$	0.7	$NP \rightarrow \textit{ears}$	0.18
$P \rightarrow VP PP$	0.3	$NP \rightarrow \textit{saw}$	0.04
$V \rightarrow \textit{with}$	1.0	$NP \rightarrow \textit{stars}$	0.18
$V \rightarrow \textit{saw}$	1.0	$NP \rightarrow \textit{telescopes}$	0.1

311



312

t_2 :



Attachment ambiguities: A key parsing decision

- The main problem in parsing is working out how to ‘attach’ various kinds of constituents – PPs, adverbial or participial phrases, coordinations, and so on
- Prepositional phrase attachment
 - *I saw the man with a telescope*
- What does *with a telescope* modify?
 - The verb *saw*?
 - The noun *man*?
- Is the problem ‘AI-complete’? Yes, but ...

Importance of lexical factors

- Words are good predictors (or even inducers) of attachment (even absent understanding):
 - The children ate the cake with a spoon.
 - The children ate the cake with frosting.
 - Moscow sent more than 100,000 soldiers into Afghanistan
 - Sydney Water breached an agreement with NSW Health
- Ford et al. (1982):
 - Ordering is jointly determined by strengths of alternative lexical forms, alternative syntactic rewrite rules, and the sequence of hypotheses in parsing

The two parse trees’ probabilities and the sentence probability

$$P(t_1) = 1.0 \times 0.1 \times 0.7 \times 1.0 \times 0.4 \times 0.18 \times 1.0 \times 1.0 \times 0.18 = 0.0009072$$

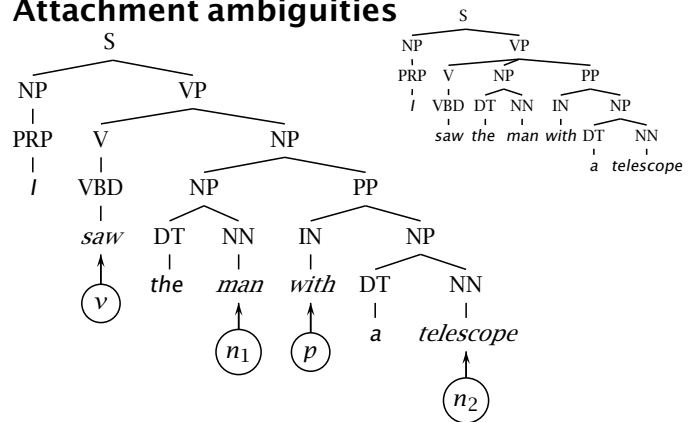
$$P(t_2) = 1.0 \times 0.1 \times 0.3 \times 0.7 \times 1.0 \times 0.18 \times 1.0 \times 1.0 \times 0.18 = 0.0006804$$

$$P(w_{15}) = P(t_1) + P(t_2) = 0.0015876$$

Attachment ambiguities (2)

- Proposed simple structural factors
 - Right association (Kimball 1973) = ‘low’ or ‘near’ attachment = ‘late closure’ (of NP) [NP → NP PP]
 - Minimal attachment (Frazier 1978) [depends on grammar] = ‘high’ or ‘distant’ attachment = ‘early closure’ (of NP) [VP → V NP PP]
- Such simple structural factors dominated in early psycholinguistics, and are still widely invoked.
- In the V NP PP context, right attachment gets it right in 55–67% of cases.
- But that means it gets it wrong in 33–45% of cases

Attachment ambiguities



Assumptions of PCFGs

1. Place invariance (like time invariance in HMM):

$$\forall k \quad P(N_{k(k+c)}^j \rightarrow \zeta) \text{ is the same}$$

2. Context-free:

$$P(N_{kl}^j \rightarrow \zeta | \text{words outside } w_k \dots w_l) = P(N_{kl}^j \rightarrow \zeta)$$

3. Ancestor-free:

$$P(N_{kl}^j \rightarrow \zeta | \text{ancestor nodes of } N_{kl}^j) = P(N_{kl}^j \rightarrow \zeta)$$

The sufficient statistics of a PCFG are thus simply counts of how often different local tree configurations occurred (= counts of which grammar rules were applied).

319

Some features of PCFGs

Reasons to use a PCFG, and some idea of their limitations:

- Partial solution for grammar ambiguity: a PCFG gives some idea of the plausibility of a sentence.
- But, in the simple case, not a very good idea, as independence assumptions are too strong (e.g., not lexicalized).
- Gives a probabilistic language model for English.
- In the simple case, a PCFG is a worse language model for English than a trigram model.
- Better for grammar induction (Gold 1967 vs. Horning 1969)
- Robustness. (Admit everything with low probability.)

321

Some features of PCFGs

- A PCFG encodes certain biases, e.g., that smaller trees are normally more probable.
- One can hope to combine the strengths of a PCFG and a trigram model.

We'll look at simple PCFGs first. They have certain inadequacies, but we'll see that most of the state-of-the-art probabilistic parsers are fundamentally PCFG models, just with various enrichments to the grammar

322

Questions for PCFGs

Just as for HMMs, there are three basic questions we wish to answer:

- Language modeling: $P(w_{1m}|G)$
- Parsing: $\arg \max_t P(t|w_{1m}, G)$
- Learning algorithm: Find G such that $P(w_{1m}|G)$ is maximized.

324

Chomsky Normal Form grammars

We'll do the case of Chomsky Normal Form grammars, which only have rules of the form:

$$N^i \rightarrow N^j N^k$$

$$N^i \rightarrow w^j$$

Any CFG can be represented by a weakly equivalent CFG in Chomsky Normal Form. It's straightforward to generalize the algorithm (recall chart parsing).

325

Probabilistic Regular Grammar:

$$N^i \rightarrow w^j N^k$$

$$N^i \rightarrow w^j$$

Start state, N^1

HMM:

$$\sum_{w_{1n}} P(w_{1n}) = 1 \quad \forall n$$

whereas in a PCFG or a PRG:

$$\sum_{w \in L} P(w) = 1$$

327

Probabilistic Regular Grammar

Consider:

$P(\text{John decided to bake a})$

High probability in HMM, low probability in a PRG or a PCFG.

Implement via sink (end) state.

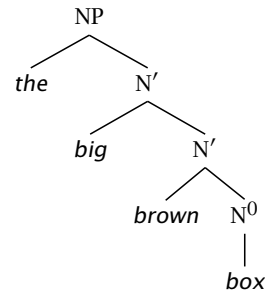
A PRG



328

Comparison of HMMs (PRGs) and PCFGs

$X: NP \rightarrow N' \rightarrow N' \rightarrow N^0 \rightarrow \text{sink}$
 $O: \text{the} \quad \text{big} \quad \text{brown} \quad \text{box}$



329

Inside and outside probabilities

This suggests: whereas for an HMM we have:

Forwards = $\alpha_i(t) = P(w_{1(t-1)}, X_t = i)$

Backwards = $\beta_i(t) = P(w_{tT} | X_t = i)$

for a PCFG we make use of Inside and Outside probabilities, defined as follows:

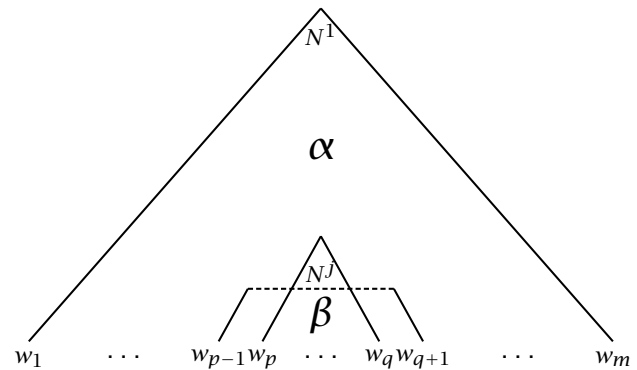
Outside = $\alpha_j(p, q) = P(w_{1(p-1)}, N_{pq}^j, w_{(q+1)m} | G)$

Inside = $\beta_j(p, q) = P(w_{pq} | N_{pq}^j, G)$

A slight generalization of dynamic Bayes Nets covers PCFG inference by the inside-outside algorithm (and-or tree of conjunctive daughters disjunctively chosen)

330

Inside and outside probabilities in PCFGs.



331

Probability of a string

Inside probability

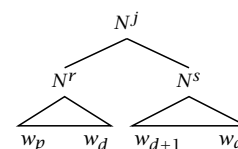
$$\begin{aligned} P(w_{1m} | G) &= P(N^1 \Rightarrow w_{1m} | G) \\ &= P(w_{1m}, N_{1m}^1, G) = \beta_1(1, m) \end{aligned}$$

Base case: We want to find $\beta_j(k, k)$ (the probability of a rule $N^j \rightarrow w_k$):

$$\begin{aligned} \beta_j(k, k) &= P(w_k | N_{kk}^j, G) \\ &= P(N^j \rightarrow w_k | G) \end{aligned}$$

332

Induction: We want to find $\beta_j(p, q)$, for $p < q$. As this is the inductive step using a Chomsky Normal Form grammar, the first rule must be of the form $N^j \rightarrow N^r N^s$, so we can proceed by induction, dividing the string in two in various places and summing the result:



These inside probabilities can be calculated bottom up.

333

For all j ,

$$\begin{aligned}
 \beta_j(p, q) &= P(w_{pq} | N_{pq}^j, G) \\
 &= \sum_{r,s} \sum_{d=p}^{q-1} P(w_{pd}, N_{pd}^r, w_{(d+1)q}, N_{(d+1)q}^s | N_{pq}^j, G) \\
 &= \sum_{r,s} \sum_{d=p}^{q-1} P(N_{pd}^r, N_{(d+1)q}^s | N_{pq}^j, G) \\
 &\quad P(w_{pd} | N_{pd}^r, N_{(d+1)q}^s, N_{(d+1)q}^s, G) \\
 &\quad P(w_{(d+1)q} | N_{pd}^r, N_{(d+1)q}^s, N_{(d+1)q}^s, w_{pd}, G) \\
 &= \sum_{r,s} \sum_{d=p}^{q-1} P(N_{pd}^r, N_{(d+1)q}^s | N_{pq}^j, G) \\
 &\quad P(w_{pd} | N_{pd}^r, G) P(w_{(d+1)q} | N_{(d+1)q}^s, G) \\
 &= \sum_{r,s} \sum_{d=p}^{q-1} P(N^j \rightarrow N^r N^s) \beta_r(p, d) \beta_s(d+1, q)
 \end{aligned}$$

334

Outside probabilities

Probability of a string: For any k , $1 \leq k \leq m$,

$$\begin{aligned}
 P(w_{1m} | G) &= \sum_j P(w_{1(k-1)}, w_k, w_{(k+1)m}, N_{kk}^j | G) \\
 &= \sum_j P(w_{1(k-1)}, N_{kk}^j, w_{(k+1)m} | G) \\
 &\quad \times P(w_k | w_{1(k-1)}, N_{kk}^j, w_{(k+1)m}, G) \\
 &= \sum_j \alpha_j(k, k) P(N^j \rightarrow w_k)
 \end{aligned}$$

Inductive (DP) calculation: One calculates the outside probabilities top down (after determining the inside probabilities).

336

Inductive Case:

$$\begin{aligned}
 \alpha_j(p, q) &= \left[\sum_{f,g} \sum_{e=q+1}^m P(w_{1(p-1)}, w_{(q+1)m}, N_{pe}^f, N_{pq}^j, N_{(q+1)e}^g) \right] \\
 &\quad + \left[\sum_{f,g} \sum_{e=1}^{p-1} P(w_{1(p-1)}, w_{(q+1)m}, N_{ef}^f, N_{(q+1)e}^g, N_{pq}^j) \right] \\
 &= \left[\sum_{f,g} \sum_{e=q+1}^m P(w_{1(p-1)}, w_{(e+1)m}, N_{pe}^f) P(N_{pq}^j, N_{(q+1)e}^g | N_{pe}^f) \right. \\
 &\quad \times P(w_{(q+1)e} | N_{(q+1)e}^g) \left. + \left[\sum_{f,g} \sum_{e=1}^{p-1} P(w_{1(e-1)}, w_{(q+1)m}, N_{ef}^f) \right. \right. \\
 &\quad \times P(N_{(q+1)e}^g | N_{pq}^j, N_{ef}^f) P(w_{(q+1)e} | N_{(q+1)e}^g) \left. \left. \right] \right] \\
 &= \left[\sum_{f,g} \sum_{e=q+1}^m \alpha_f(p, e) P(N^f \rightarrow N^j) \beta_g(q+1, e) \right] \\
 &\quad + \left[\sum_{f,g} \sum_{e=1}^{p-1} \alpha_f(e, q) P(N^f \rightarrow N^j) \beta_g(e, p-1) \right]
 \end{aligned}$$

340

Calculation of inside probabilities (CKY algorithm)

	1	2	3	4	5
1	$\beta_{NP} = 0.1$		$\beta_S = 0.0126$		$\beta_S = 0.0015876$
2		$\beta_{NP} = 0.04$ $\beta_V = 1.0$	$\beta_{VP} = 0.126$		$\beta_{VP} = 0.015876$
3			$\beta_{NP} = 0.18$		$\beta_{NP} = 0.01296$
4				$\beta_P = 1.0$	$\beta_{PP} = 0.18$
5					$\beta_{NP} = 0.18$
	astronomers	saw	stars	with	ears

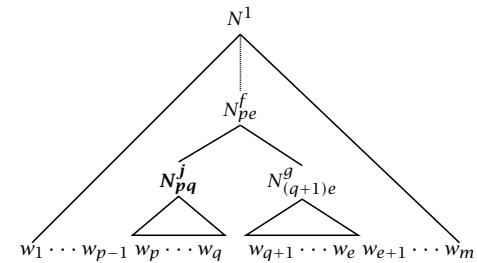
335

Outside probabilities

Base Case:

$$\begin{aligned}
 \alpha_1(1, m) &= 1 \\
 \alpha_j(1, m) &= 0, \text{ for } j \neq 1
 \end{aligned}$$

Inductive Case: it's either a left or right branch - we will solve over both possibilities and calculate using outside and inside probabilities



338

Overall probability of a node existing

As with a HMM, we can form a product of the inside and outside probabilities. This time:

$$\begin{aligned}
 \alpha_j(p, q) \beta_j(p, q) &= P(w_{1(p-1)}, N_{pq}^j, w_{(q+1)m} | G) P(w_{pq} | N_{pq}^j, G) \\
 &= P(w_{1m}, N_{pq}^j | G)
 \end{aligned}$$

Therefore,

$$p(w_{1m}, N_{pq} | G) = \sum_j \alpha_j(p, q) \beta_j(p, q)$$

Just in the cases of the root node and the preterminals, we know there will always be some such constituent.

341

Finding the most likely parse (Viterbi algorithm)

Like inside algorithm, but find maximum rather than sum
Record which rule gave this maximum

$\delta_i(p, q)$ = the highest inside probability parse of a subtree N_{pq}^i

1. Initialization: $\delta_i(p, p) = P(N^i \rightarrow w_p)$

2. Induction

$$\delta_i(p, q) = \max_{\substack{1 \leq j, k \leq n \\ p \leq r < q}} P(N^i \rightarrow N^j N^k) \delta_j(p, r) \delta_k(r + 1, q)$$

3. Store backtrace

$$\psi_i(p, q) = \arg \max_{(j, k, r)} P(N^i \rightarrow N^j N^k) \delta_j(p, r) \delta_k(r + 1, q)$$

4. From start symbol N^1 , most likely parse t is:

t begins with $\psi_1(1, m)$. $P(\hat{t}) = \delta_1(1, m)$

342

Learning PCFGs (1)

- We would like to calculate how often each rule is used:

$$\hat{P}(N^j \rightarrow \zeta) = \frac{C(N^j \rightarrow \zeta)}{\sum_{\gamma} C(N^j \rightarrow \gamma)}$$

- If we have labeled data, we count and find out
- Relative frequency again gives maximum likelihood probability estimates
- This is the motivation for building *Treebanks* of hand-parsed sentences

351

Calculation of Viterbi probabilities (CKY algorithm)

	1	2	3	4	5
1	$\delta_{NP} = 0.1$		$\delta_S = 0.0126$		$\delta_S = 0.0009072$
2		$\delta_{NP} = 0.04$ $\delta_V = 1.0$	$\delta_{VP} = 0.126$		$\delta_{VP} = 0.009072$
3			$\delta_{NP} = 0.18$		$\delta_{NP} = 0.01296$
4				$\delta_P = 1.0$	$\delta_{PP} = 0.18$
5					$\delta_{NP} = 0.18$
	astronomers	saw	stars	with	ears

343

Learning PCFGs (2): the Inside-Outside algorithm (Baker 1979)

- Otherwise we work iteratively from expectations of current model.
- We construct an EM training algorithm, as for HMMs
- For each sentence, at each iteration, we work out expectation of how often each rule is used using inside and outside probabilities
- We assume sentences are independent and sum expectations over parses of each
- We re-estimate rules based on these 'counts'

352