

Statistical NLP Spring 2009



Lecture 15: PCFGs

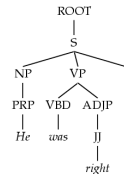
Dan Klein – UC Berkeley



Treebank PCFGs

[Charniak 96]

- Use PCFGs for broad coverage parsing
- Can take a grammar right off the trees (doesn't work well):

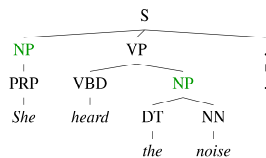


- ROOT → S 1
- S → NP VP 1
- NP → PRP 1
- VP → VBD ADJP 1
-

Model	F1
Baseline	72.0



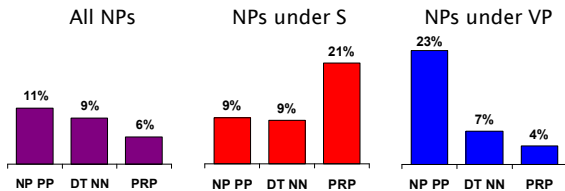
Conditional Independence?



- Not every NP expansion can fill every NP slot
 - A grammar with symbols like "NP" won't be context-free
 - Statistically, conditional independence too strong

Non-Independence

- Independence assumptions are often too strong.

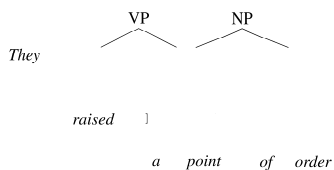


- Example: the expansion of an NP is highly dependent on the parent of the NP (i.e., subjects vs. objects).
- Also: the subject and object expansions are correlated!

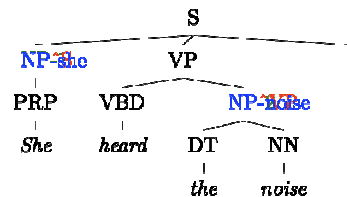


Grammar Refinement

- Example: PP attachment

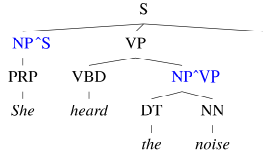


Grammar Refinement



- Structure Annotation [Johnson '98, Klein&Manning '03]
- Lexicalization [Collins '99, Charniak '00]
- Latent Variables [Matsuzaki et al. '05, Petrov et al. '06]

The Game of Designing a Grammar



- Annotation refines base treebank symbols to improve statistical fit of the grammar
 - Structural annotation

Typical Experimental Setup

- Corpus: Penn Treebank, WSJ

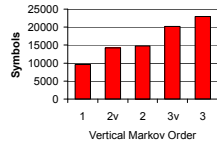
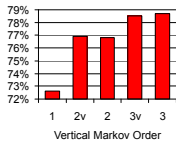
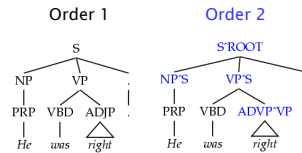


Training: sections 02-21
 Development: section 22 (here, first 20 files)
 Test: section 23

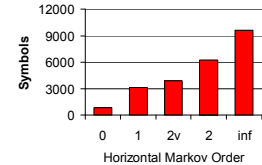
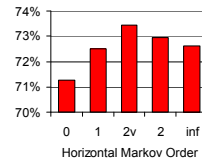
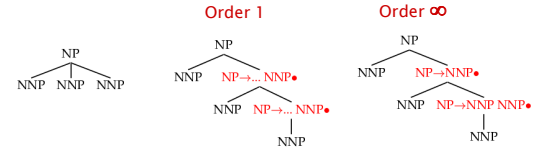
- Accuracy – F1: harmonic mean of per-node labeled precision and recall.
- Here: also size – number of symbols in grammar.
 - Passive / complete symbols: NP, NP^S
 - Active / incomplete symbols: NP → NP CC •

Vertical Markovization

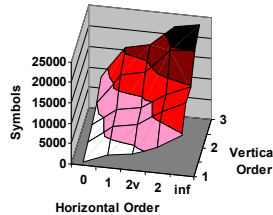
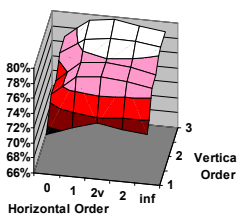
- Vertical Markov order: rewrites depend on past k ancestor nodes. (cf. parent annotation)



Horizontal Markovization



Vertical and Horizontal

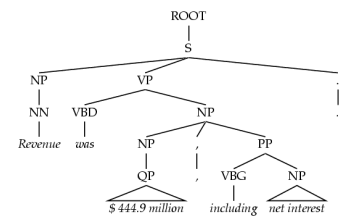


- Examples:
 - Raw treebank: $v=1, h=\infty$
 - Johnson 98: $v=2, h=\infty$
 - Collins 99: $v=2, h=2$
 - Best F1: $v=3, h=2v$

Model	F1	Size
Base: $v=h=2v$	77.8	7.5K

Unary Splits

- Problem: unary rewrites used to transmute categories so a high-probability rule can be used.
- Solution: Mark unary rewrite sites with -U



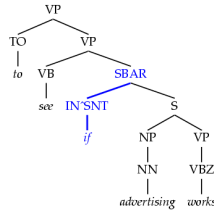
Annotation	F1	Size
Base	77.8	7.5K
UNARY	78.3	8.0K

Tag Splits

- Problem: Treebank tags are too coarse.

- Example: Sentential, PP, and other prepositions are all marked IN.

- Partial Solution:
 - Subdivide the IN tag.



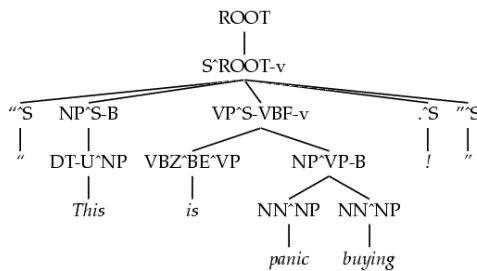
Annotation	F1	Size
Previous	78.3	8.0K
SPLIT-IN	80.3	8.1K

Other Tag Splits

- UNARY-DT: mark demonstratives as DT^U ("the X" vs. "those")
- UNARY-RB: mark phrasal adverbs as RB^U ("quickly" vs. "very")
- TAG-PA: mark tags with non-canonical parents ("not" is an RB^VP)
- SPLIT-AUX: mark auxiliary verbs with -AUX [cf. Charniak 97]
- SPLIT-CC: separate "but" and "&" from other conjunctions
- SPLIT-%: "%" gets its own tag.

F1	Size
80.4	8.1K
80.5	8.1K
81.2	8.5K
81.6	9.0K
81.7	9.1K
81.8	9.3K

A Fully Annotated (Unlex) Tree

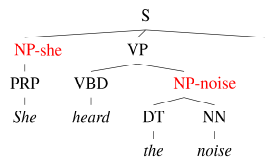


Some Test Set Results

Parser	LP	LR	F1	CB	0 CB
Magerman 95	84.9	84.6	84.7	1.26	56.6
Collins 96	86.3	85.8	86.0	1.14	59.9
Unlexicalized	86.9	85.7	86.3	1.10	60.3
Charniak 97	87.4	87.5	87.4	1.00	62.1
Collins 99	88.7	88.6	88.6	0.90	67.1

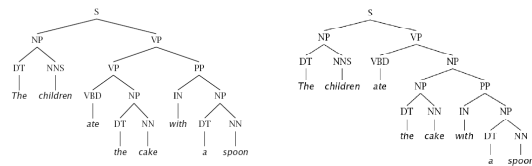
- Beats "first generation" lexicalized parsers.
- Lots of room to improve – more complex models next.

The Game of Designing a Grammar



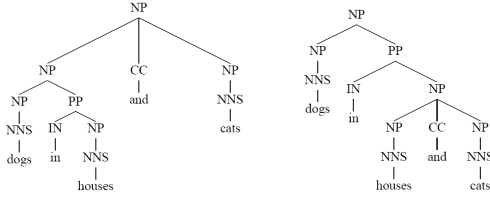
- Annotation refines base treebank symbols to improve statistical fit of the grammar
 - Structural annotation [Johnson '98, Klein and Manning 03]
 - Head lexicalization [Collins '99, Charniak '00]

Problems with PCFGs



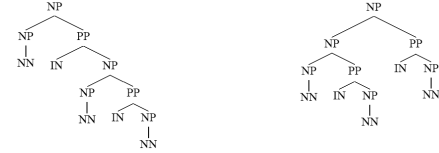
- If we do no annotation, these trees differ only in one rule:
 - VP → VP PP
 - NP → NP PP
- Parse will go one way or the other, regardless of words
- We addressed this in one way with unlexicalized grammars (how?)
- Lexicalization allows us to be sensitive to specific words

Problems with PCFGs



- What's different between basic PCFG scores here?
- What (lexical) correlations need to be scored?

Problems with PCFGs

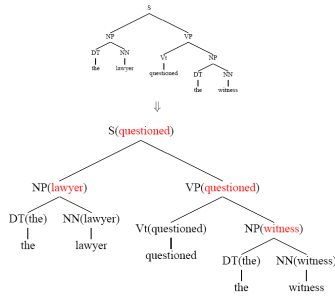


president of a company in Africa

- Another example of PCFG indifference
 - Left structure far more common
 - How to model this?
 - Really structural: "chicken with potatoes with gravy"
 - Lexical parsers model this effect, but not by virtue of being lexical

Lexicalized Trees

- Add "headwords" to each phrasal node
 - Syntactic vs. semantic heads
 - Headship not in (most) treebanks
 - Usually use head rules, e.g.:
 - NP:
 - Take leftmost NP
 - Take rightmost N*
 - Take rightmost JJ
 - Take right child
 - VP:
 - Take leftmost VB*
 - Take leftmost VP
 - Take left child



Lexicalized PCFGs?

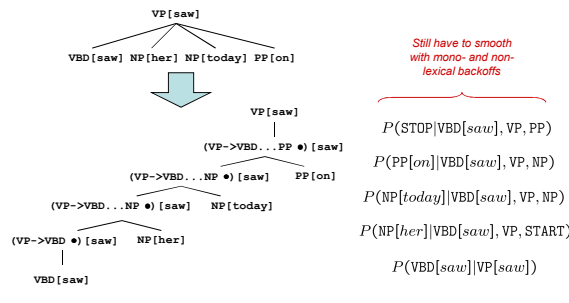
- Problem: we now have to estimate probabilities like

$$VP(\text{saw}) \rightarrow VBD(\text{saw}) NP-C(\text{her}) NP(\text{today})$$
- Never going to get these atomically off of a treebank
- Solution: break up derivation into smaller steps



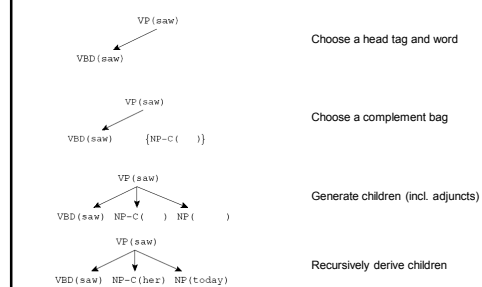
Lexical Derivation Steps

- Derivation of a local tree [simplified Charniak 97]



Lexical Derivation Steps

- Another derivation of a local tree [Collins 99]



Naïve Lexicalized Parsing

- Can, in principle, use CKY on lexicalized PCFGs
 - $O(Rn^3)$ time and $O(Sn^2)$ memory
 - But $R = rV^2$ and $S = sV$
 - Result is completely impractical (why?)
 - Memory: 10K rules * 50K words * (40 words)² * 8 bytes ≈ 6TB
- Can modify CKY to exploit lexical sparsity
 - Lexicalized symbols are a base grammar symbol and a pointer into the input sentence, not any arbitrary word
 - Result: $O(m^5)$ time, $O(sn^3)$
 - Memory: 10K rules * (40 words)³ * 8 bytes ≈ 5GB

Lexicalized CKY

(VP->VBD ... NP) [saw]

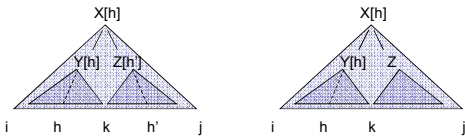
(VP->VBD) [saw] NP [hez]

```

bestScore(X, i, j, h)
  if (j = i+1)
    return tagScore(X, s[i])
  else
    return
      max_{k, h'} max_{h'} score(X[h]->Y[h] Z[h']) *
        bestScore(Y, i, k, h) *
        bestScore(Z, k, j, h')
      max_{k, h'} score(X[h]->Y[h'] Z[h]) *
        bestScore(Y, i, k, h') *
        bestScore(Z, k, j, h)
    
```

Quartic Parsing

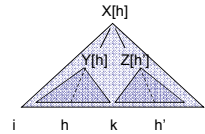
- Turns out, you can do (a little) better [Eisner 99]



- Gives an $O(n^4)$ algorithm
- Still prohibitive in practice if not pruned

Pruning with Beams

- The Collins parser prunes with per-cell beams [Collins 99]
 - Essentially, run the $O(n^5)$ CKY
 - Remember only a few hypotheses for each span $\langle i, j \rangle$.
 - If we keep K hypotheses at each span, then we do at most $O(nK^2)$ work per span (why?)
 - Keeps things more or less cubic
- Also: certain spans are forbidden entirely on the basis of punctuation (crucial for speed)

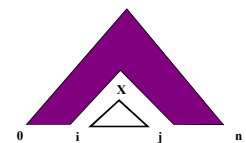


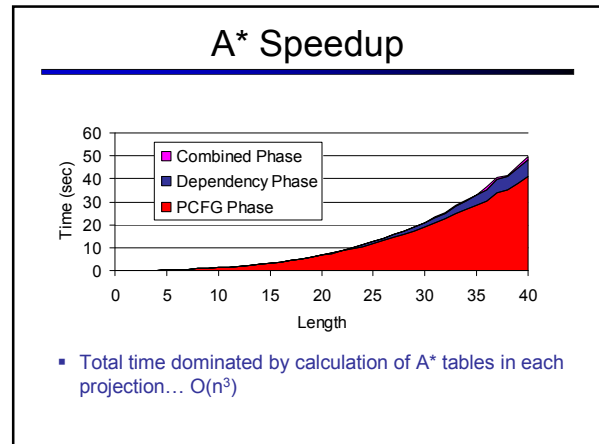
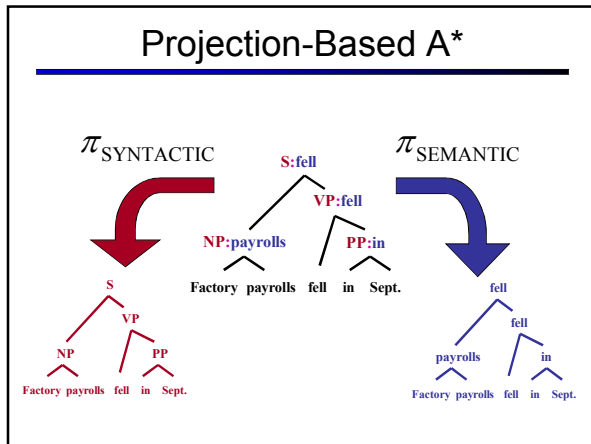
Pruning with a PCFG

- The Charniak parser prunes using a two-pass approach [Charniak 97+]
 - First, parse with the base grammar
 - For each $X: [i, j]$ calculate $P(X|i, j, s)$
 - This isn't trivial, and there are clever speed ups
 - Second, do the full $O(n^5)$ CKY
 - Skip any $X: [i, j]$ which had low (say, < 0.0001) posterior
 - Avoids almost all work in the second phase!
- Charniak et al 06: can use more passes
- Petrov et al 07: can use many more passes

Pruning with A*

- You can also speed up the search without sacrificing optimality
- For agenda-based parsers:
 - Can select which items to process first
 - Can do with any "figure of merit" [Charniak 98]
 - If your figure-of-merit is a valid A* heuristic, no loss of optimality [Klein and Manning 03]





- ### Results
- Some results
 - Collins 99 – 88.6 F1 (generative lexical)
 - Charniak and Johnson 05 – 89.7 / 91.3 F1 (generative lexical / reranked)
 - Petrov et al 06 – 90.7 F1 (generative unlexical)
 - McClosky et al 06 – 92.1 F1 (gen + rerank + self-train)
 - However
 - Bilexical counts rarely make a difference (why?)
 - Gildea 01 – Removing bilexical counts costs < 0.5 F1
 - Bilexical vs. monolexical vs. smart smoothing

The Game of Designing a Grammar

S
NP-1 VP
PRP VBD DT NN
She heard the noise

- Annotation refines base treebank symbols to improve statistical fit of the grammar
 - Structural annotation
 - Head lexicalization
 - Automatic clustering?

Latent Variable Grammars

University of California Berkeley

Parse Tree T
Sentence w

Derivations $t: T$

Parameters θ

Grammar G

$S_1 \rightarrow NP_0 VP_0 ?$
 $S_0 \rightarrow NP_1 VP_0 ?$
 $S_0 \rightarrow NP_0 VP_1 ?$
 $S_0 \rightarrow NP_1 VP_1 ?$
 $S_1 \rightarrow NP_2 VP_0 ?$
 $S_1 \rightarrow NP_2 VP_1 ?$
 $S_2 \rightarrow NP_1 VP_1 ?$
 \dots
 $NP_0 \rightarrow PRP_0 ?$
 $NP_0 \rightarrow NP_1 ?$
 \dots
 Leadcom
 $PRP_1 \rightarrow She ?$
 $PRP_1 \rightarrow She ?$
 \dots
 $VBD_0 \rightarrow was ?$
 $VBD_1 \rightarrow was ?$
 $VBD_2 \rightarrow was ?$
 \dots

Learning Latent Annotations

EM algorithm:

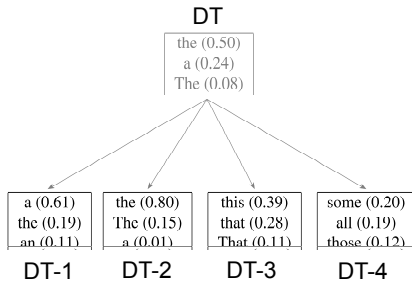
- Brackets are known
- Base categories are known
- Only induce subcategories

Forward

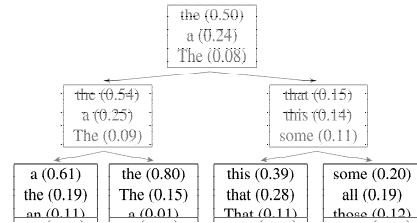
Backward

Just like Forward-Backward for HMMs.

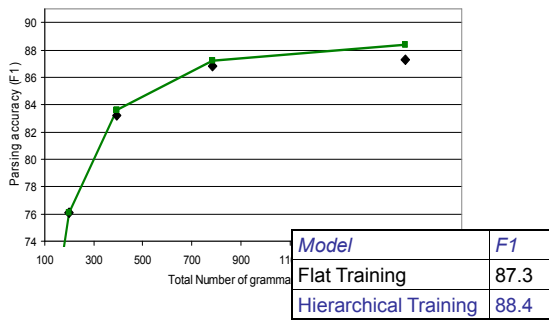
Refinement of the DT tag



Hierarchical refinement

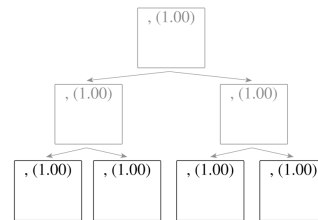


Hierarchical Estimation Results



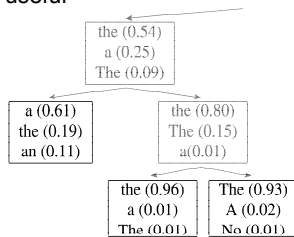
Refinement of the , tag

- Splitting all categories equally is wasteful:

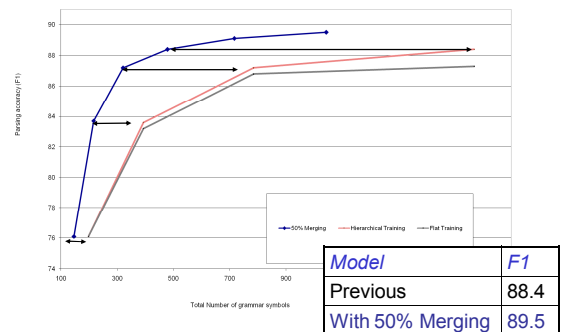


Adaptive Splitting

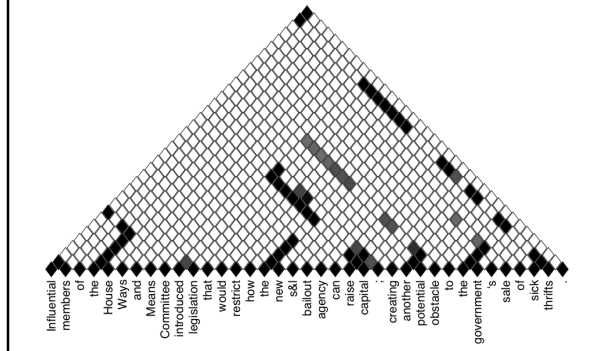
- Want to split complex categories more
- Idea: split everything, roll back splits which were least useful



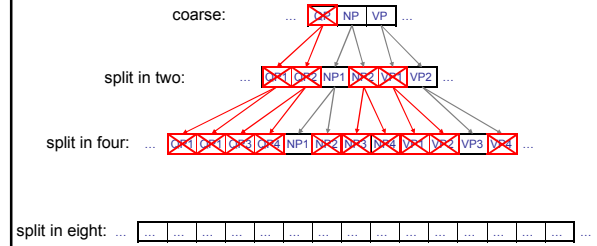
Adaptive Splitting Results



Bracket Posteriors



Hierarchical Pruning



Final Results (Accuracy)

		≤ 40 words F1	all F1
ENG	Charniak&Johnson '05 (generative)	90.1	89.6
	Split / Merge	90.6	90.1
GER	Dubey '05	76.3	-
	Split / Merge	80.8	80.1
CHN	Chiang et al. '02	80.0	76.6
	Split / Merge	86.3	83.4

Still higher numbers from reranking / self-training methods