

Statistical NLP Spring 2011



Lecture 16: Parsing II

Dan Klein – UC Berkeley

Classical NLP: Parsing

- Write symbolic or logical rules:

Grammar (CFG)		Lexicon
ROOT → S	NP → NP PP	NN → interest
S → NP VP	VP → VBP NP	NNS → raises
NP → DT NN	VP → VBP NP PP	VBP → interest
NP → NN NNS	PP → IN NP	VBZ → raises
		...

- Use deduction systems to prove parses from words
 - Minimal grammar on "Fed raises" sentence: 36 parses
 - Simple 10-rule grammar: 592 parses
 - Real-size grammar: many millions of parses
- This scaled very badly, didn't yield broad-coverage tools

A Recursive Parser

```
bestScore(X,i,j,s)
  if (j = i+1)
    return tagScore(X,s[i])
  else
    return max score(X->YZ) *
             bestScore(Y,i,k) *
             bestScore(Z,k,j)
```

- Will this parser work?
- Why or why not?
- Memory requirements?

A Memoized Parser

- One small change:

```
bestScore(X,i,j,s)
  if (scores[X][i][j] == null)
    if (j = i+1)
      score = tagScore(X,s[i])
    else
      score = max score(X->YZ) *
              bestScore(Y,i,k) *
              bestScore(Z,k,j)
    scores[X][i][j] = score
  return scores[X][i][j]
```

A Bottom-Up Parser (CKY)

- Can also organize things bottom-up

```
bestScore(s)
  for (i : [0,n-1])
    for (X : tags[s[i]])
      score[X][i][i+1] =
        tagScore(X,s[i])
  for (diff : [2,n])
    for (i : [0,n-diff])
      j = i + diff
      for (X->YZ : rule)
        for (k : [i+1, j-1])
          score[X][i][j] = max score[X][i][k],
                             score[X->YZ) *
                             score[Y][i][k] *
                             score[Z][k][j]
```

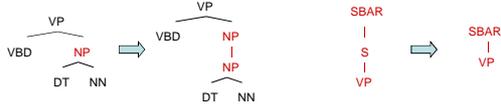
Unary Rules

- Unary rules?

```
bestScore(X,i,j,s)
  if (j = i+1)
    return tagScore(X,s[i])
  else
    return max max score(X->YZ) *
              bestScore(Y,i,k) *
              bestScore(Z,k,j)
             max score(X->Y) *
             bestScore(Y,i,j)
```

CNF + Unary Closure

- We need unaries to be non-cyclic
 - Can address by pre-calculating the *unary closure*
 - Rather than having zero or more unaries, always have exactly one



- Alternate unary and binary layers
- Reconstruct unary chains afterwards

Alternating Layers

```
bestScoreB(X,i,j,s)
    return max max score(X->YZ) *
                bestScoreU(Y,i,k) *
                bestScoreU(Z,k,j)
```

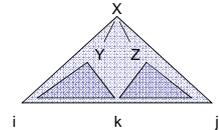
```
bestScoreU(X,i,j,s)
    if (j = i+1)
        return tagScore(X,s[i])
    else
        return max max score(X->Y) *
                    bestScoreB(Y,i,j)
```

Memory

- How much memory does this require?
 - Have to store the score cache
 - Cache size: |symbols|*n² doubles
 - For the plain treebank grammar:
 - X ~ 20K, n = 40, double ~ 8 bytes = ~ 256MB
 - Big, but workable.
- Pruning: Beams
 - score[X][i][j] can get too large (when?)
 - Can keep beams (truncated maps score[i][j]) which only store the best few scores for the span [i,j]
- Pruning: Coarse-to-Fine
 - Use a smaller grammar to rule out most X[i,j]
 - Much more on this later...

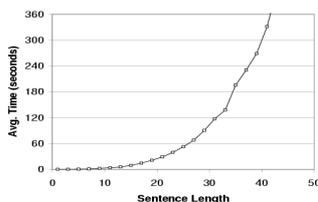
Time: Theory

- How much time will it take to parse?
 - For each diff (<= n)
 - For each i (<= n)
 - For each rule X → Y Z
 - For each split point k
 - Do constant work
- Total time: |rules|*n³
- Something like 5 sec for an unoptimized parse of a 20-word sentences



Time: Practice

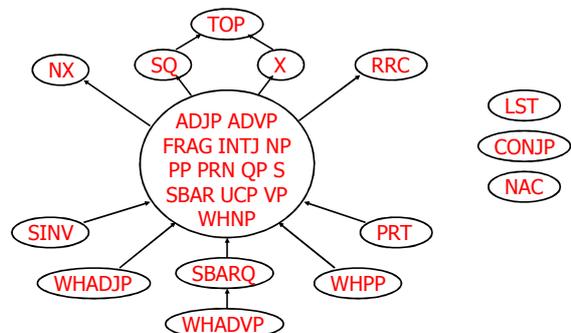
- Parsing with the vanilla treebank grammar:



~ 20K Rules
(not an optimized parser!)
Observed exponent:
3.6

- Why's it worse in practice?
 - Longer sentences "unlock" more of the grammar
 - All kinds of systems issues don't scale

Same-Span Reachability

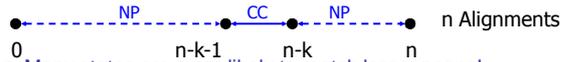


Rule State Reachability

Example: NP CC •



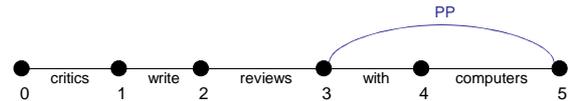
Example: NP CC NP •



- Many states are more likely to match larger spans!

Agenda-Based Parsing

- Agenda-based parsing is like graph search (but over a hypergraph)
- Concepts:
 - Numbering: we number fenceposts between words
 - "Edges" or items: spans with labels, e.g. PP[3,5], represent the sets of trees over those words rooted at that label (cf. search states)
 - A chart: records edges we've expanded (cf. closed set)
 - An agenda: a queue which holds edges (cf. a fringe or open set)



Word Items

- Building an item for the first time is called discovery. Items go into the agenda on discovery.
- To initialize, we discover all word items (with score 1.0).

AGENDA

critics[0,1], write[1,2], reviews[2,3], with[3,4], computers[4,5]

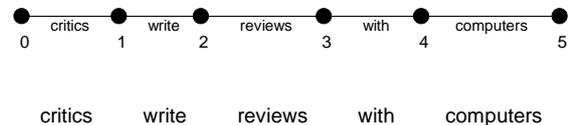
CHART [EMPTY]



Unary Projection

- When we pop a word item, the lexicon tells us the tag item successors (and scores) which go on the agenda

critics[0,1] write[1,2] reviews[2,3] with[3,4] computers[4,5]
 NNS[0,1] VBP[1,2] NNS[2,3] IN[3,4] NNS[4,5]



Item Successors

- When we pop items off of the agenda:
 - Graph successors: unary projections (NNS → critics, NP → NNS)

$Y[i,j]$ with $X \rightarrow Y$ forms $X[i,j]$

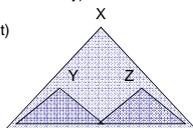
- Hypergraph successors: combine with items already in our chart

$Y[i,j]$ and $Z[j,k]$ with $X \rightarrow YZ$ form $X[i,k]$

- Enqueue / promote resulting items (if not in chart already)
- Record backtraces as appropriate
- Stick the popped edge in the chart (closed set)

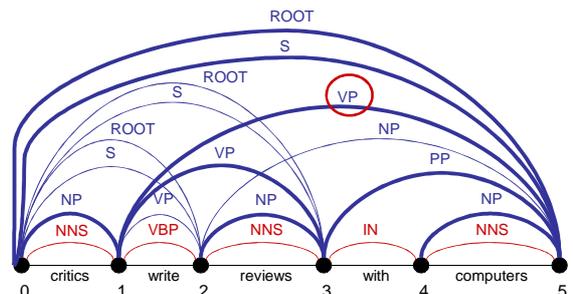
- Queries a chart must support:

- Is edge $X[i,j]$ in the chart? (What score?)
- What edges with label Y end at position j ?
- What edges with label Z start at position i ?



An Example

NNS[0,1] VBP[1,2] NNS[2,3] IN[3,4] NNS[4,5] NP[0,1] VP[1,2] NP[2,3] NP[4,5] S[0,2]
 VP[1,3] PP[3,5] ROOT[0,2] S[0,3] VP[1,5] NP[2,5] ROOT[0,3] S[0,5] ROOT[0,5]



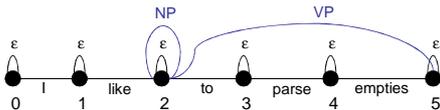
Empty Elements

- Sometimes we want to posit nodes in a parse tree that don't contain any pronounced words:

I want you to parse this sentence

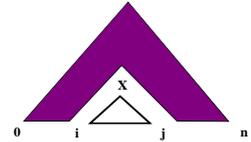
I want [] to parse this sentence

- These are easy to add to a chart parser!
 - For each position i , add the "word" edge $\epsilon:[i,i]$
 - Add rules like $NP \rightarrow \epsilon$ to the grammar
 - That's it!



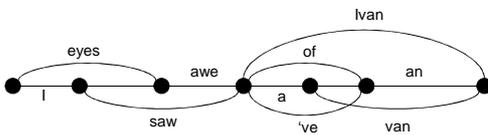
UCS / A*

- With weighted edges, order matters
 - Must expand optimal parse from bottom up (subparses first)
 - CKY does this by processing smaller spans before larger ones
 - UCS pops items off the agenda in order of decreasing $Viterbi$ score
 - A* search also well defined
- You can also speed up the search without sacrificing optimality
 - 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]



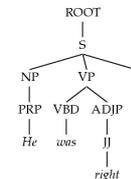
(Speech) Lattices

- There was nothing magical about words spanning exactly one position.
- When working with speech, we generally don't know how many words there are, or where they break.
- We can represent the possibilities as a lattice and parse these just as easily.



Treebank PCFGs [Charniak 96]

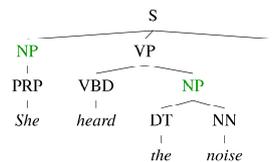
- Use PCFGs for broad coverage parsing
- Can take a grammar right off the trees (doesn't work well):



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

Model	F1
Baseline	72.0

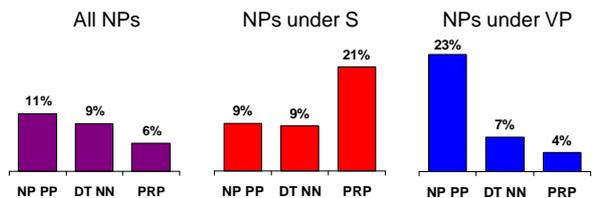
Conditional Independence?



- 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.

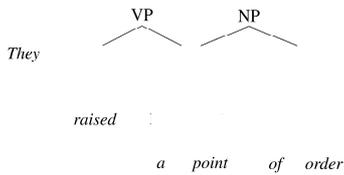


- Example: the expansion of an NP is highly dependent on the parent of the NP (i.e., subjects vs. objects).
- Also: the subject and object expansions are correlated!

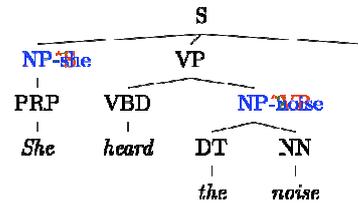


Grammar Refinement

- Example: PP attachment

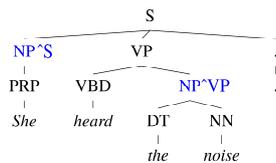


Grammar Refinement



- 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

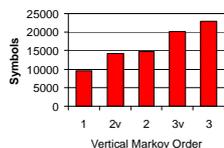
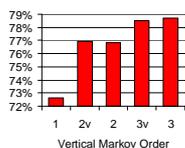
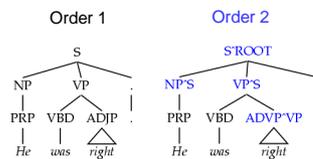
- Corpus: Penn Treebank, WSJ



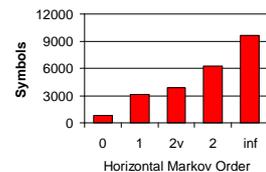
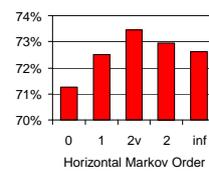
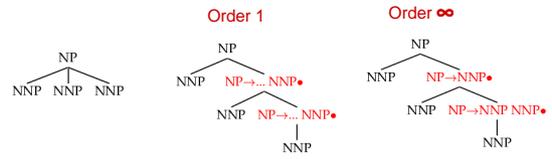
- Accuracy – F1: harmonic mean of per-node labeled precision and recall.
- Here: also size – number of symbols in grammar.
 - Passive / complete symbols: NP, NP^S
 - Active / incomplete symbols: NP → NP CC •

Vertical Markovization

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

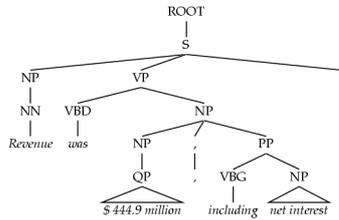


Horizontal Markovization



Unary Splits

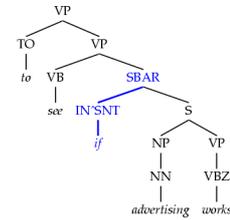
- Problem: unary rewrites used to transmute categories so a high-probability rule can be used.
- Solution: Mark unary rewrite sites with -U



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.



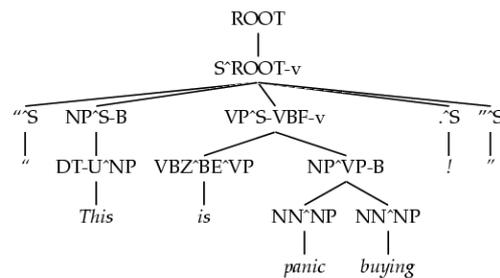
Annotation	F1	Size
Previous	78.3	8.0K
SPLIT-IN	80.3	8.1K

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^U 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
UNARY-DT	80.4	8.1K
UNARY-RB	80.5	8.1K
TAG-PA	81.2	8.5K
SPLIT-AUX	81.6	9.0K
SPLIT-CC	81.7	9.1K
SPLIT-%	81.8	9.3K

A Fully Annotated (Unlex) Tree

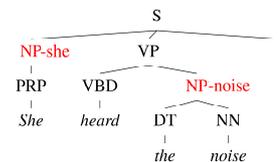


Some Test Set Results

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

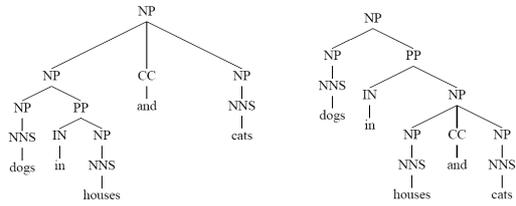
- Beats "first generation" lexicalized parsers.
- Lots of room to improve – more complex models next.

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]

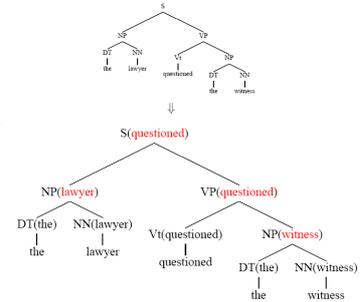
Problems with PCFGs



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

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

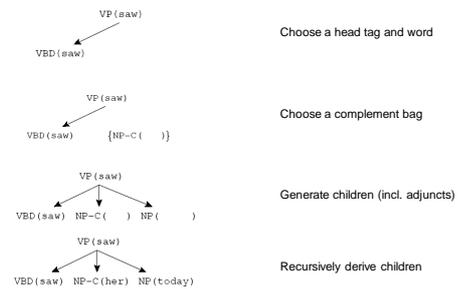
$VP(\text{saw}) \rightarrow VBD(\text{saw}) NP-C(\text{her}) NP(\text{today})$

- Never going to get these atomically off of a treebank
- Solution: break up derivation into smaller steps



Lexical Derivation Steps

- A derivation of a local tree [Collins 99]

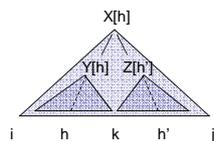


Lexicalized CKY

$(VP \rightarrow VBD \dots NP) [saw]$
 $(VP \rightarrow VBD) [saw] \quad NP[her]$

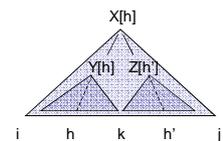
```

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



Pruning with Beams

- The Collins parser prunes with per-cell beams [Collins 99]
 - Essentially, run the $O(n^3)$ 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



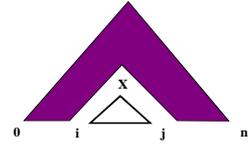
- Also: certain spans are forbidden entirely on the basis of punctuation (crucial for speed)

Pruning with a PCFG

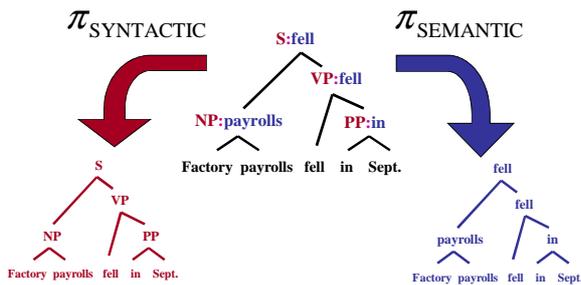
- The Charniak parser prunes using a two-pass approach [Charniak 97+]
 - First, parse with the base grammar
 - For each $X:[i,j]$ calculate $P(X|i,j,s)$
 - This isn't trivial, and there are clever speed ups
 - Second, do the full $O(n^5)$ CKY
 - Skip any $X:[i,j]$ which had low (say, < 0.0001) posterior
 - 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*

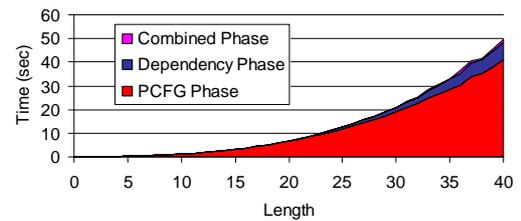
- You can also speed up the search without sacrificing optimality
- For agenda-based parsers:
 - Can select which items to process first
 - Can do with any "figure of merit" [Charniak 98]
 - If your figure-of-merit is a valid A* heuristic, no loss of optimality [Klein and Manning 03]



Projection-Based A*



A* Speedup



- Total time dominated by calculation of A* tables in each projection... $O(n^3)$

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