
Speech and Language Processing: Statistical Parsing

Chapter 14

1

Statistical Parsing

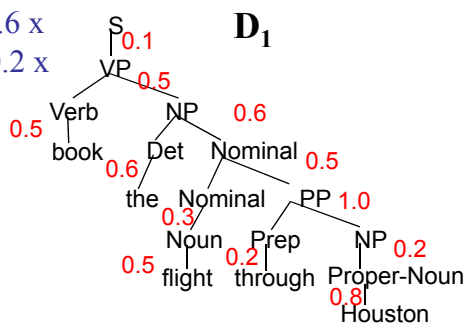
- Statistical parsing uses a probabilistic model of syntax in order to assign probabilities to each parse tree.
- Provides principled approach to resolving syntactic ambiguity.
- Allows supervised learning of parsers from tree-banks of parse trees provided by human linguists.

2

Sentence Probability

- Assume productions for each node are chosen independently.
- Probability of derivation is the product of the probabilities of its productions.

$$\begin{aligned}
 P(D_1) &= 0.1 \times 0.5 \times 0.5 \times 0.6 \times 0.6 \times \\
 &\quad 0.5 \times 0.3 \times 1.0 \times 0.2 \times 0.2 \times \\
 &\quad 0.5 \times 0.8 \\
 &= 0.0000216
 \end{aligned}$$



5

Other Parses?

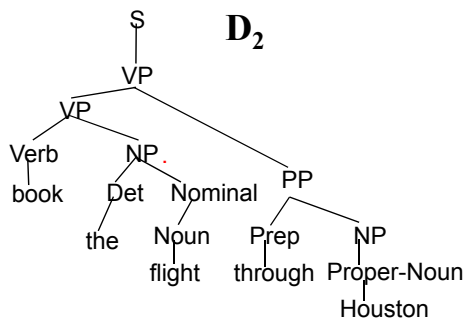
- book the flight through Houston

6

Syntactic Disambiguation

- Resolve ambiguity by picking most probable parse tree.

$P(D_2) =$

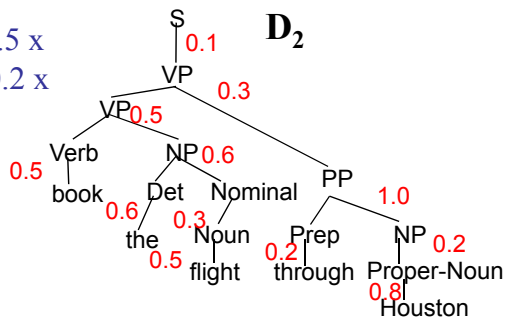


7

Syntactic Disambiguation

- Resolve ambiguity by picking most probable parse tree.

$$\begin{aligned}
 P(D_2) &= 0.1 \times 0.3 \times 0.5 \times 0.6 \times 0.5 \times \\
 &\quad 0.6 \times 0.3 \times 1.0 \times 0.5 \times 0.2 \times \\
 &\quad 0.2 \times 0.8 \\
 &= 0.00001296
 \end{aligned}$$

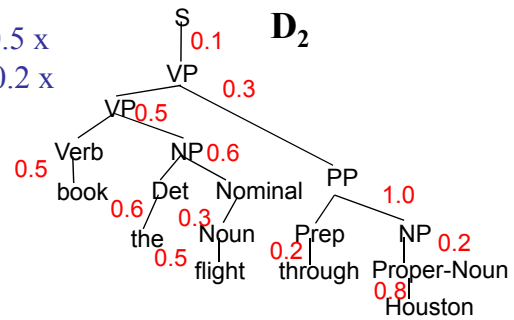


8

Syntactic Disambiguation

- Resolve ambiguity by picking most probable parse tree.

$$\begin{aligned}
 P(D_2) &= 0.1 \times 0.3 \times 0.5 \times 0.6 \times 0.5 \times \\
 &\quad 0.6 \times 0.3 \times 1.0 \times 0.5 \times 0.2 \times \\
 &\quad 0.2 \times 0.8 \\
 &= 0.00001296
 \end{aligned}$$



9

Disambiguation Result?

10

Sentence Probability

- Probability of a sentence is the sum of the probabilities of all of its derivations.

$$\begin{aligned} P(\text{"book the flight through Houston"}) &= \\ P(D_1) + P(D_2) &= 0.0000216 + 0.00001296 \\ &= 0.00003456 \end{aligned}$$

11

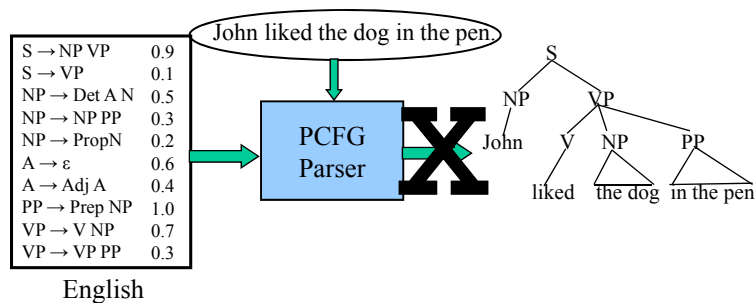
Three Useful PCFG Tasks

- **Observation likelihood:** To classify and order sentences.
- **Most likely derivation:** To determine the most likely parse tree for a sentence.
- **Maximum likelihood training:** To train a PCFG to fit empirical training data.

12

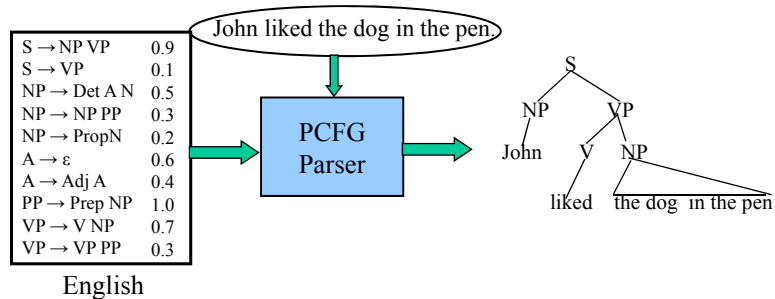
PCFG: Most Likely Derivation

- There is an analog to the Viterbi algorithm to efficiently determine the most probable derivation (parse tree) for a sentence.



PCFG: Most Likely Derivation

- There is an analog to the Viterbi algorithm to efficiently determine the most probable derivation (parse tree) for a sentence.



Probabilistic CKY

- CKY can be modified for PCFG parsing by including in each cell a probability for each non-terminal.
- Cell $[i,j]$ must retain the *most probable* derivation of each constituent (non-terminal) covering words $i + 1$ through j together with its associated probability.
- When transforming the grammar to CNF, must set production probabilities to preserve the probability of derivations.

Probabilistic Grammar Conversion

Original Grammar	Probability	Chomsky Normal Form	Probability
$S \rightarrow NP VP$	0.8	$S \rightarrow NP VP$	0.8
$S \rightarrow Aux NP VP$	0.1	$S \rightarrow X1 VP$	0.1
		$X1 \rightarrow Aux NP$	1.0
$S \rightarrow VP$	0.1	$S \rightarrow \text{book} \mid \text{include} \mid \text{prefer}$	
		0.01 0.004 0.006	
		$S \rightarrow \text{Verb NP}$	0.05
		$S \rightarrow VP PP$	0.03
$NP \rightarrow \text{Pronoun}$	0.2	$NP \rightarrow I \mid he \mid she \mid me$	
		0.1 0.02 0.02 0.06	
$NP \rightarrow \text{Proper-Noun}$	0.2	$NP \rightarrow \text{Houston} \mid \text{NWA}$	
		0.16 .04	
$NP \rightarrow \text{Det Nominal}$	0.6	$NP \rightarrow \text{Det Nominal}$	0.6
$Nominal \rightarrow \text{Noun}$	0.3	$Nominal \rightarrow \text{book} \mid \text{flight} \mid \text{meal} \mid \text{money}$	
		0.03 0.15 0.06 0.06	
$Nominal \rightarrow \text{Nominal Noun}$	0.2	$Nominal \rightarrow \text{Nominal Noun}$	0.2
$Nominal \rightarrow \text{Nominal PP}$	0.5	$Nominal \rightarrow \text{Nominal PP}$	0.5
$VP \rightarrow \text{Verb}$	0.2	$VP \rightarrow \text{book} \mid \text{include} \mid \text{prefer}$	
		0.1 0.04 0.06	
$VP \rightarrow \text{Verb NP}$	0.5	$VP \rightarrow \text{Verb NP}$	0.5
$VP \rightarrow VP PP$	0.3	$VP \rightarrow VP PP$	0.3
$PP \rightarrow \text{Prep NP}$	1.0	$PP \rightarrow \text{Prep NP}$	1.0

Probabilistic CKY Parser

Book the flight through Houston

S :.01, VP:.1, Verb:.5 Nominal:.03 Noun:.1	None			
	Det:.6	NP:.6*.6*.15 =.054		
		Nominal:.15 Noun:.5		

17

Probabilistic CKY Parser

Book the flight through Houston

S :.01, VP:.1, Verb:.5 Nominal:.03 Noun:.1	None	VP:.5*.5*.054 =.0135		
	Det:.6	NP:.6*.6*.15 =.054		
		Nominal:.15 Noun:.5		

18

Probabilistic CKY Parser

Book the flight through Houston

S :.01, VP:.1, Verb:.5 Nominal:.03 Noun:.1	None	S:.05*.5*.054 =.00135 VP:.5*.5*.054 =.0135	
	Det:.6	NP:.6*.6*.15 =.054	
		Nominal:.15 Noun:.5	

19

Probabilistic CKY Parser

Book the flight through Houston

S :.01, VP:.1, Verb:.5 Nominal:.03 Noun:.1	None	S:.05*.5*.054 =.00135 VP:.5*.5*.054 =.0135	None
	Det:.6	NP:.6*.6*.15 =.054	None
		Nominal:.15 Noun:.5	None
			Prep:.2

20

Probabilistic CKY Parser

Book the flight through Houston

S :.01, VP:.1, Verb:.5 Nominal:.03 Noun:.1	None	S:.05*.5*.054 =.00135 VP:.5*.5*.054 =.0135	None	
	Det:.6	NP:.6*.6*.15 =.054	None	
		Nominal:.15 Noun:.5	None	
			Prep:.2 ←	PP:1.0*.2*.16 =.032
				NP:.16 PropNoun:.8

21

Probabilistic CKY Parser

Book the flight through Houston

S :.01, VP:.1, Verb:.5 Nominal:.03 Noun:.1	None	S:.05*.5*.054 =.00135 VP:.5*.5*.054 =.0135	None	
	Det:.6	NP:.6*.6*.15 =.054	None	
		Nominal:.15 Noun:.5	← None	Nominal: .5*.15*.032 =.0024
			Prep:.2	PP:1.0*.2*.16 =.032
				NP:.16 PropNoun:.8

22

Probabilistic CKY Parser

Book the flight through Houston

S :.01, VP:.1, Verb:.5 Nominal:.03 Noun:.1	None	S:.05*.5*.054 =.00135 VP:.5*.5*.054 =.0135	None	
	Det:.6 ←	NP:.6*.6*.15 =.054	None	NP:.6*.6* .0024 =.000864
		Nominal:.15 Noun:.5	None	Nominal: .5*.15*.032 =.0024
			Prep:.2	PP:1.0*.2*.16 =.032
				NP:.16 PropNoun:.8

23

Probabilistic CKY Parser

Book the flight through Houston

S :.01, VP:.1, Verb: .5 Nominal:.03 Noun:.1	None	S:.05*.5*.054 =.00135 VP:.5*.5*.054 =.0135	None	S:.05*.5* .000864 =.000216
	Det:.6	NP:.6*.6*.15 =.054	None	NP:.6*.6* .0024 =.000864
		Nominal:.15 Noun:.5	None	Nominal: .5*.15*.032 =.0024
			Prep:.2	PP:1.0*.2*.16 =.032
				NP:.16 PropNoun:.8

24

Probabilistic CKY Parser

Book the flight through Houston

S :.01, VP:.1, Verb:.5 Nominal:.03 Noun:.1	None	S:.05*.5*.054 =.00135 VP:.5*.5*.054 =.0135	None	S:.03*.0135* .032 =.0001296 S:.0000216
	Det:.6	NP:.6*.6*.15 =.054	None	NP:.6*.6* .0024 =.000864
		Nominal:.15 Noun:.5	None	Nominal: .5*.15*.032 =.0024
			Prep:.2	PP:1.0*.2*.16 =.032
				NP:.16 PropNoun:.8

25

Probabilistic CKY Parser

Book the flight through Houston

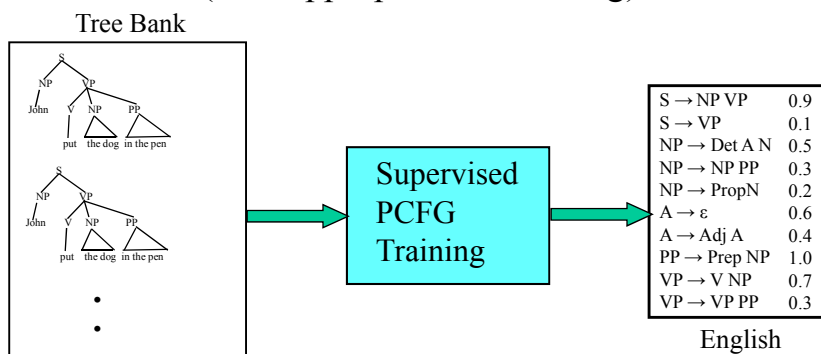
S :.01, VP:.1, Verb: .5 Nominal:.03 Noun:.1	None	S: .05*.5*.054 =.00135 VP: .5*.5*.054 =.0135	None	S:.0000216
	Det: .6	NP: .6*.6*.15 =.054	None	NP: .6*.6* .0024 =.000864
		Nominal: .15 Noun: .5	None	Nominal: .5*.15*.032 =.0024
			Prep: .2	PP:1.0*.2*.16 =.032
				NP:.16 PropNoun:.8

Pick most probable parse, i.e. take max to combine probabilities of multiple derivations of each constituent in each cell.

26

PCFG: Supervised Training

- If parse trees are provided for training sentences, a grammar and its parameters can be estimated directly from counts accumulated from the **tree-bank** (with appropriate smoothing).



27

Estimating Production Probabilities

- Set of production rules can be taken directly from the set of rewrites in the treebank.
- Parameters can be directly estimated from frequency counts in the treebank.

$$P(\alpha \rightarrow \beta | \alpha) = \frac{\text{count}(\alpha \rightarrow \beta)}{\sum_{\gamma} \text{count}(\alpha \rightarrow \gamma)} = \frac{\text{count}(\alpha \rightarrow \beta)}{\text{count}(\alpha)}$$

28

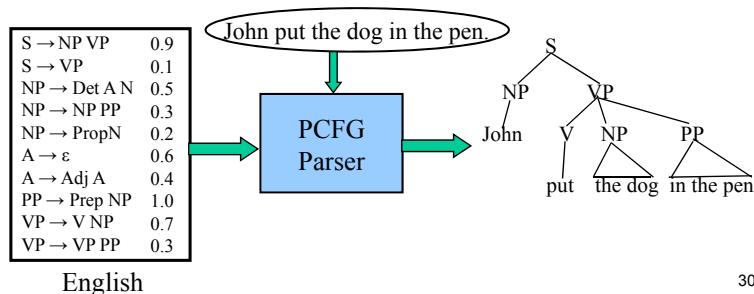
Vanilla PCFG Limitations

- Since probabilities of productions do not rely on specific words or concepts, only general structural disambiguation is possible (e.g. prefer to attach PPs to Nominals).
- Consequently, vanilla PCFGs cannot resolve syntactic ambiguities that require semantics to resolve, e.g. ate with fork vs. meatballs.
- In order to work well, PCFGs must be **lexicalized**, i.e. productions must be specialized to specific words by including their head-word in their LHS non-terminals (e.g. VP-ate).

29

Example of Importance of Lexicalization

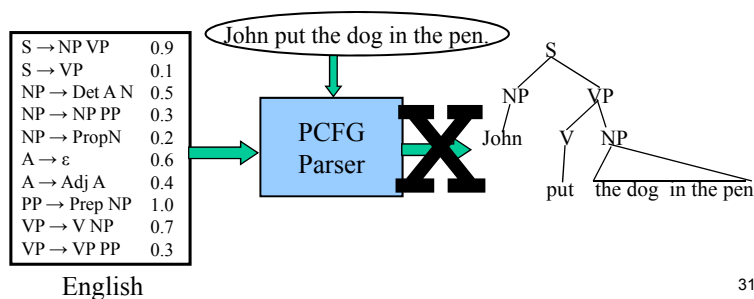
- A general preference for attaching PPs to NPs rather than VPs can be learned by a vanilla PCFG.
- But the desired preference can depend on specific words.



30

Example of Importance of Lexicalization

- A general preference for attaching PPs to NPs rather than VPs can be learned by a vanilla PCFG.
- But the desired preference can depend on specific words.



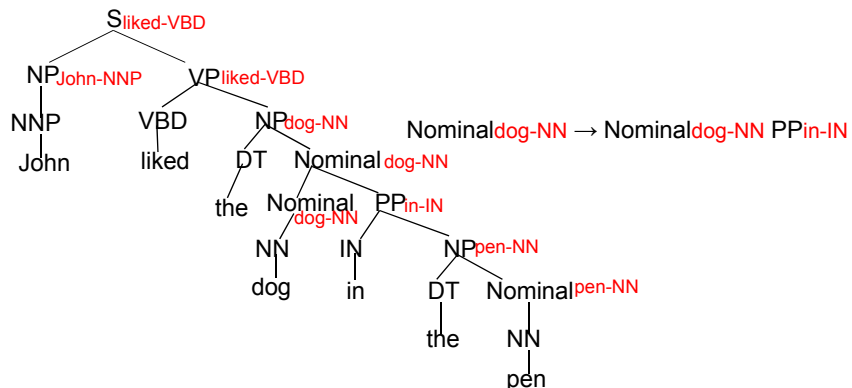
31

Head Words

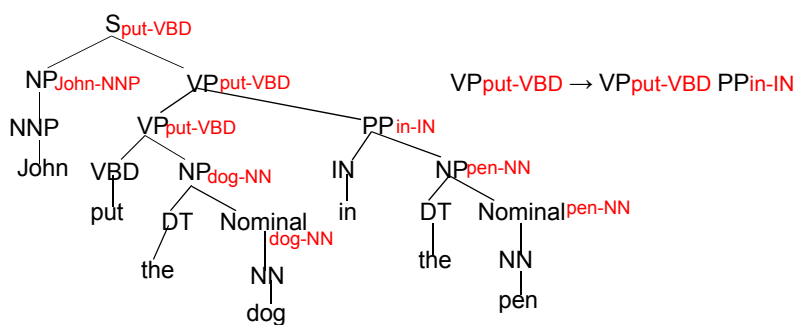
- Syntactic phrases usually have a word in them that is most “central” to the phrase.
- Linguists have defined the concept of a lexical **head** of a phrase.
- Simple rules can identify the head of any phrase by percolating head words up the parse tree.
 - Head of a VP is the main verb
 - Head of an NP is the main noun
 - Head of a PP is the preposition
 - Head of a sentence is the head of its VP

Lexicalized Productions

- Specialized productions can be generated by including the head word and its POS of each non-terminal as part of that non-terminal's symbol.



Lexicalized Productions



Parameterizing Lexicalized Productions

- Accurately estimating parameters on such a large number of very specialized productions could require enormous amounts of treebank data.
- Need some way of estimating parameters for lexicalized productions that makes reasonable independence assumptions so that accurate probabilities for very specific rules can be learned.

Collins' Parser

- Collins' (1999) parser assumes a simple generative model of lexicalized productions.
- Models productions based on context to the left and the right of the head daughter.
 - $LHS \rightarrow L_n L_{n-1} \dots L_1 H R_1 \dots R_{m-1} R_m$
- First generate the head (H) and then repeatedly generate left (L_i) and right (R_i) context symbols until the symbol STOP is generated.

Missed Context Dependence

- Another problem with CFGs is that which production is used to expand a non-terminal is independent of its context.
- However, this independence is frequently violated for normal grammars.
 - NPs that are subjects are more likely to be pronouns than NPs that are objects.

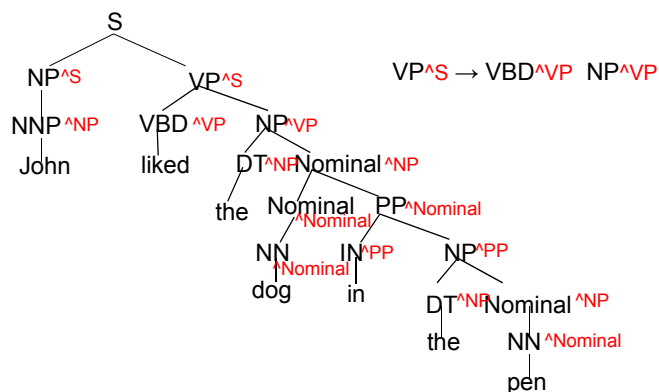
39

Splitting Non-Terminals

- To provide more contextual information, non-terminals can be split into multiple new non-terminals based on their parent in the parse tree using **parent annotation**.
 - A subject NP becomes NP^S since its parent node is an S.
 - An object NP becomes NP^{VP} since its parent node is a VP

40

Parent Annotation Example



41

Split and Merge

- Non-terminal splitting greatly increases the size of the grammar and the number of parameters that need to be learned from limited training data.
- Best approach is to only split non-terminals when it improves the accuracy of the grammar.
- May also help to merge some non-terminals to remove some un-helpful distinctions and learn more accurate parameters for the merged productions.
- Method: Heuristically search for a combination of splits and merges that produces a grammar that maximizes the likelihood of the training treebank.

42

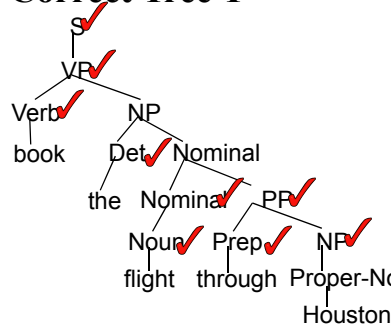
Parsing Evaluation Metrics

- PARSEVAL metrics measure the fraction of the constituents that match between the computed and human parse trees. If P is the system's parse tree and T is the human parse tree (the "gold standard"):
 - **Recall** = (# correct constituents in P) / (# constituents in T)
 - **Precision** = (# correct constituents in P) / (# constituents in P)
- **Labeled Precision** and **labeled recall** require getting the non-terminal label on the constituent node correct to count as correct.
- F_1 is the harmonic mean of precision and recall.

43

Computing Evaluation Metrics

Correct Tree T



Constituents: 12

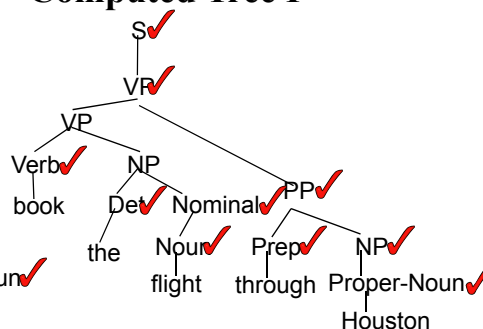
Correct Constituents: 10

Recall = 10/12 = 83.3%

Precision = 10/12 = 83.3%

 $F_1 = 83.3\%$

Computed Tree P



Constituents: 12

Recall = 10/12 = 83.3%

Precision = 10/12 = 83.3%

 $F_1 = 83.3\%$

Treebank Results

- Results of current state-of-the-art systems on the English Penn WSJ treebank are slightly greater than 90% labeled precision and recall.

45

Discriminative Parse Reranking

- Motivation: Even when the top-ranked parse not correct, frequently the correct parse is one of those ranked highly by a statistical parser.
- Use a discriminative classifier that is trained to select the best parse from the N-best parses produced by the original parser.
- Reranker can exploit global features of the entire parse whereas a PCFG is restricted to making decisions based on local info.

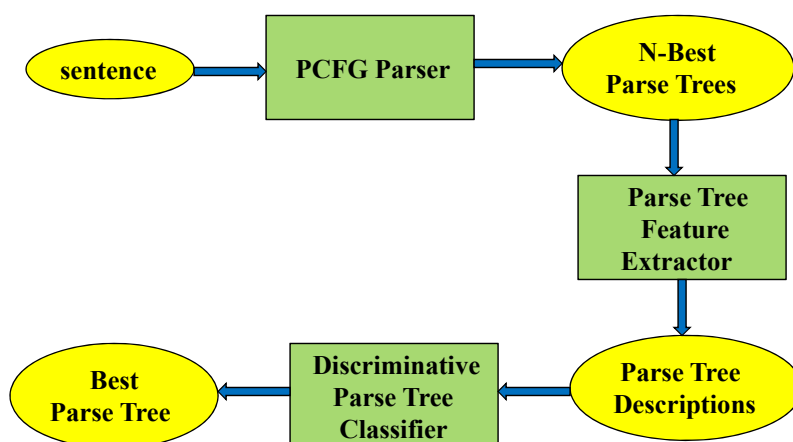
46

2-Stage Reranking Approach

- Adapt the PCFG parser to produce an *N-best list* of the most probable parses in addition to the most-likely one.
- Extract from each of these parses, a set of global features that help determine if it is a good parse tree.
- Train a discriminative classifier (e.g. logistic regression) using the best parse in each N-best list as positive and others as negative.

47

Parse Reranking



48

Sample Parse Tree Features

- Probability of the parse from the PCFG.
- The number of parallel conjuncts.
 - “the bird in the tree and the squirrel on the ground”
 - “the bird and the squirrel in the tree”
- The degree to which the parse tree is right branching.
 - English parses tend to be right branching (cf. parse of “Book the flight through Houston”)

49

Evaluation of Reranking

- Reranking is limited by *oracle accuracy*, i.e. the accuracy that results when an omniscient oracle picks the best parse from the N-best list.
- Typical current oracle accuracy is around $F_1=97\%$
- Reranking can generally improve test accuracy of current PCFG models a percentage point or two.

50

Human Parsing

- Computational parsers can be used to predict human reading time as measured by tracking the time taken to read each word in a sentence.
- Psycholinguistic studies show that words that are more probable given the preceding lexical and syntactic context are read faster.

51

Garden Path Sentences

- People are confused by sentences that seem to have a particular syntactic structure but then suddenly violate this structure, so the listener is “lead down the garden path”.
 - The horse raced past the barn fell.
 - vs. The horse raced past the barn broke his leg.
 - The complex houses married students.
 - The old man the sea.
 - While Anna dressed the baby spit up on the bed.

52

Unification Grammars (Ch 15)

- In order to handle agreement issues more effectively, each constituent has a list of features such as number, person, gender, etc. which may or not be specified for a given constituent.
- In order for two constituents to combine to form a larger constituent, their features must *unify*, i.e. consistently combine into a merged set of features.
- Expressive grammars and parsers (e.g. HPSG – head driven phrase structure grammar) have been developed using this approach and have been partially integrated with modern statistical models of disambiguation.

53

Mildly Context-Sensitive Grammars

- Some grammatical formalisms provide a degree of context-sensitivity that helps capture aspects of NL syntax that are not easily handled by CFGs.
- Combinatory Categorical Grammar (CCG) consists of:
 - **Categorial Lexicon** that associates a syntactic and semantic category with each word.
 - **Combinatory Rules** that define how categories combine to form other categories.

54

Statistical Parsing Conclusions

- Statistical models such as PCFGs allow for probabilistic resolution of ambiguities.
- PCFGs can be easily learned from treebanks.
- Lexicalization and non-terminal splitting are required to effectively resolve many ambiguities.
- Current statistical parsers are quite accurate but not yet at the level of human-expert agreement.

55