

# Statistical parsing <sup>2</sup>

---

- Using statistical models for
  - Ambiguity resolution
    - E.g. Pp-attachment
  - Guiding parsing
  - Get the most likely parse
- Grammatical induction from corpora
- Goal: Parsing of non restricted texts with a reasonable level of accuracy (>90%) and efficiency.
- Requirements:
  - Corpora tagged (with POS): Brown, LOB, Clic-Talp
  - Corpora analyzed: Penn treebank, Susanne, Ancora

# Statistical parsing <sup>3</sup>

---

- Lexical approaches
  - Context free: unigram
  - Context dependent: N-gram, HMM
- Syntactic approaches
  - SCFG
- Hybrid approaches
  - Stochastic lexicalized Tags
- Computing the most likely (most probable) parse
  - Viterbi
- Parameter learning
  - Supervised
    - Tagged corpora
  - Non supervised
    - Baum-Welch (Fw-Bw) para HMM
    - Inside-Outside para SCFG

# SCFG<sub>1</sub>

---

- Associate a probability to each rule
- Associate a probability to each lexical entry
- Frequent restriction CNF:
  - Binary rules  $A_p \rightarrow A_q A_r$  matrix  $B_{pqr}$
  - Unary rules  $A_p \rightarrow b_m$  matrix  $U_{pm}$

# SCFG <sub>2</sub>

---

S	⇒	NP	VP	1.0
VP	⇒	Vi		0.4
VP	⇒	Vt	NP	0.4
VP	⇒	VP	PP	0.2
NP	⇒	DT	NN	0.3
NP	⇒	NP	PP	0.7
PP	⇒	P	NP	1.0

Vi	⇒	sleeps	1.0
Vt	⇒	saw	1.0
NN	⇒	man	0.7
NN	⇒	woman	0.2
NN	⇒	telescope	0.1
DT	⇒	the	1.0
IN	⇒	with	0.5
IN	⇒	in	0.5

- Probability of a tree  $t$  with rules

$$\alpha_1 \rightarrow \beta_1, \alpha_2 \rightarrow \beta_2, \dots, \alpha_n \rightarrow \beta_n$$

is

$$p(t) = \prod_{i=1}^n q(\alpha_i \rightarrow \beta_i)$$

where  $q(\alpha \rightarrow \beta)$  is the probability for rule  $\alpha \rightarrow \beta$ .

# SCFG <sub>3</sub>

---

1. A context-free grammar  $G = (N, \Sigma, S, R)$ .

2. A parameter

$$q(\alpha \rightarrow \beta)$$

for each rule  $\alpha \rightarrow \beta \in R$ . The parameter  $q(\alpha \rightarrow \beta)$  can be interpreted as the conditional probability of choosing rule  $\alpha \rightarrow \beta$  in a left-most derivation, given that the non-terminal being expanded is  $\alpha$ . For any  $X \in N$ , we have the constraint

$$\sum_{\alpha \rightarrow \beta \in R: \alpha = X} q(\alpha \rightarrow \beta) = 1$$

In addition we have  $q(\alpha \rightarrow \beta) \geq 0$  for any  $\alpha \rightarrow \beta \in R$ .

Given a parse-tree  $t \in \mathcal{T}_G$  containing rules  $\alpha_1 \rightarrow \beta_1, \alpha_2 \rightarrow \beta_2, \dots, \alpha_n \rightarrow \beta_n$ , the probability of  $t$  under the PCFG is

$$p(t) = \prod_{i=1}^n q(\alpha_i \rightarrow \beta_i)$$

# SCFG<sub>4</sub>

---

- Assigns a probability to each *left-most derivation*, or parse-tree, allowed by the underlying CFG
- Say we have a sentence  $s$ , set of derivations for that sentence is  $\mathcal{T}(s)$ . Then a PCFG assigns a probability  $p(t)$  to each member of  $\mathcal{T}(s)$ . i.e., *we now have a ranking in order of probability*.
- The most likely parse tree for a sentence  $s$  is

$$\arg \max_{t \in \mathcal{T}(s)} p(t)$$

# SCFG<sub>5</sub>

---

- ▶ **Learning.** How to obtain a PCFG from a treebank?
- ▶ **Inference.** Given a PCFG and a sentence  $s$ . Denote by  $\mathcal{T}(s)$  the set of derivations that yield  $s$ .
  - ▶ How to compute the best parse for  $s$ ?

$$\operatorname{argmax}_{t \in \mathcal{T}(s)} p(t)$$

- ▶ How to compute the probability of  $s$ ?

$$\sum_{t \in \mathcal{T}(s)} p(t)$$

# SCFG<sub>6</sub>

---

## Pros and cons of SCFG

- Some idea of the probability of a parse
  - But not very good.
- CFG cannot be learned without negative examples, SCFG can
- SCFGs provide a LM for a language
- In practice SCFG provide a worse LM than a 3-gram
  - $P([\text{N} [\text{N} \text{toy}] [\text{N} [\text{N} \text{coffee}] [\text{N} \text{grinder}]]]) = P([\text{N} [\text{N} [\text{N} \text{cat}] [\text{N} \text{food}]]] [\text{N} \text{tin}])$
  - $P(\text{NP} \rightarrow \text{Pro})$  is  $>$  in Subj position than in Obj position.



# SCFG <sup>7</sup>

---

- Robust
- Possibility of combining SCFG with 3-grams
- SCFG assign a lot of probability mass to short sentences (a small tree is more probable than a big one)
- Parameter estimation (probabilities)
- Problem of sparseness
- Volume

## SCFG <sup>8</sup>

---

- Association to the rule. Information about the point of application of the rule in the derivation tree is lost.
- Low frequency constructions are penalized
- Probability of a derivation. Contextual independence is assumed
- Possibility of relax conditional independence:
  - Sensitivity to structure
  - Lexicalization

# SCFG <sup>9</sup>

---

- Sensitivity to structure
    - Node expansion depends of its position in the tree
- |         | <u>Pronoun</u> | <u>Lexical</u> |
|---------|----------------|----------------|
| Subject | 91%            | 9%             |
| Object  | 34%            | 66%            |
- Enrichment of a node with information of its ancestors
    - sNP is different of <sup>v</sup>PNP
  - Pronouns as arguments of ditransitive verbs
    - *I gave Charlie the book*
    - *I gave the book to Charlie*
    - *I gave you the book*
    - ? *I gave the book to you*

- (Head) Lexicalization
  - *put* takes both an NP and a VP
    - *Sue put* [ *the book* ]<sub>NP</sub> [ *on the table* ]<sub>PP</sub>
    - \* *Sue put* [ *the book* ]<sub>NP</sub>
    - \* *Sue put* [ *on the table* ]<sub>PP</sub>
  - *like* usually takes an NP and not a PP
    - *Sue likes* [ *the book* ]<sub>NP</sub>
    - \* *Sue likes* [ *on the table* ]<sub>PP</sub>

# SCFG <sup>11</sup>

---

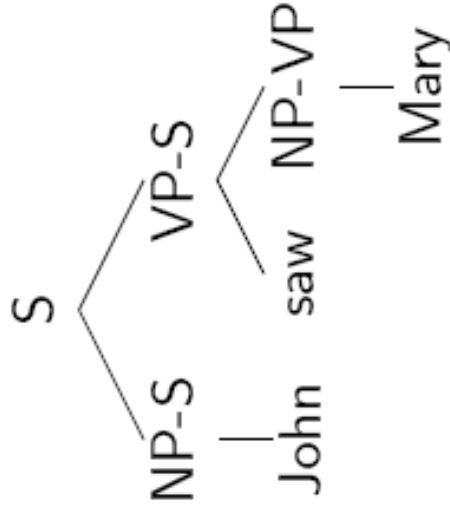
<i>Local Tree</i>	<i>Come</i>	<i>Take</i>	<i>Think</i>	<i>Want</i>
VP-> V	9.5%	2.6%	4.6%	5.7%
VP-> V NP	1.1%	32.1%	0.2%	13.9%
VP-> V PP	34.5%	3.1%	7.1%	0.3%
VP- V SBAR	6.6%	0.3%	73.0%	0.2%
VP-> V S	2.2%	1.3%	4.8%	70.8%
VP->V NP S	0.1%	5.7%	0.0%	0.3%
VP->V PRT NP	0.3%	5.8%	0.0%	0.0%
VP->V PRT PP	6.1%	1.5%	0.2%	0.0%

# SCFG<sub>12</sub>

---

## PCFGs with Parent Annotations

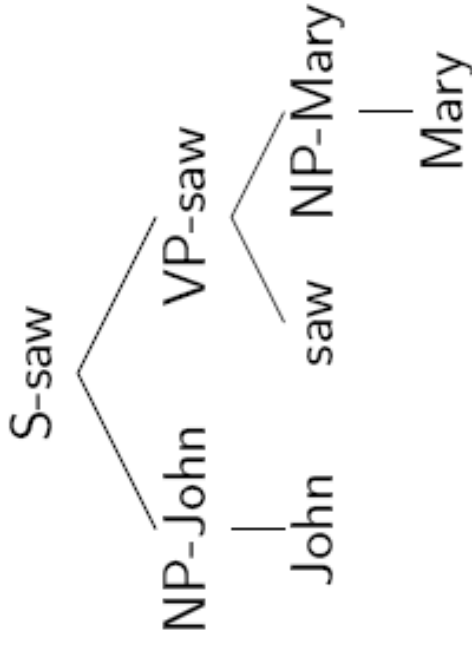
(Johnson, 1999)



$$P(\text{Tree}) = q(S \rightarrow \text{NP-S VP-S} \mid S) \times q(\text{NP-S} \rightarrow \text{John} \mid \text{NP-S}) \times \dots$$

## Lexicalized PCFGs

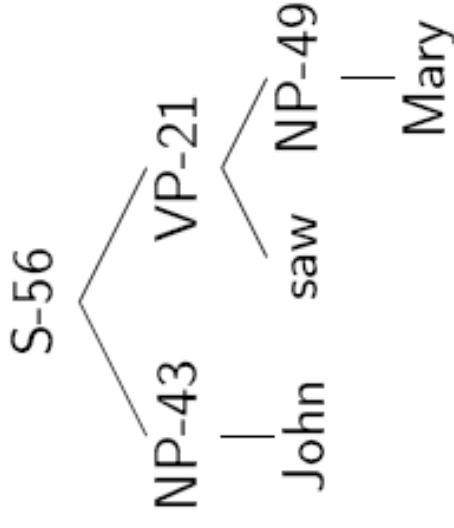
(Collins, 1999)



$$\begin{aligned}
 P(\text{Tree}) &= P(\text{S-saw} \rightarrow \text{NP-John VP-saw} \mid \text{S-saw}) \times \\
 &\quad P(\text{NP-John} \rightarrow \text{John} \mid \text{NP-John}) \times \\
 &\quad \dots
 \end{aligned}$$

## PCFGs with Latent Variables

(e.g., Petrov and Klein, 2007)



- ▶ Each non-terminals (e.g., S) is split into a number of new non-terminals (e.g., S-1, S-2, ..., S-128)
- ▶ Latent annotations learned using EM



# Parsing SCFG 1

---

- Two models
- **Conditional/Discriminative model :**
- The probability of a parse tree is directly estimated

$$\hat{t} = \underset{t}{\operatorname{argmax}} P(t | s, G) \text{ con } \sum_t P(t | s, G) = 1$$

- Probabilities are conditioned on a concrete sentence.
- No sentence distribution probabilities is assumed

- **Generative/joined model :**
  - Assigns probabilities to all the trees generated by the grammar. Probabilities are, thus, for language  $L$ :

$$\sum_{\{t: \text{yield}(t) \in L\}} P(t) = 1$$

$$\hat{t} = \underset{t}{\operatorname{argmax}} P(t, s | G) = \underset{t}{\operatorname{argmax}} P(t | s, G) / P(s | G)$$

# Parsing SCFG <sub>3</sub>

---

- CFG
- SCFG
- For each rule of G,  $(A \rightarrow \alpha) \in P_G$  we should be able to define a probability  $P(A \rightarrow \alpha)$

$$\sum_{(A \rightarrow \alpha) \in P_G} P(A \rightarrow \alpha) = 1$$

- Probability of a tree

$$P(\psi) = \prod_{(A \rightarrow \alpha) \in P_G} P(A \rightarrow \alpha)^{f(A \rightarrow \alpha; \psi)}$$

- $P(t)$  -- Probability of a tree  $t$  (product of probabilities of the rules generating it).
- $P(w_{In})$  -- Probability of a sentence is the sum of the probabilities of all the valid parse trees of the sentence

$$\begin{aligned} P(w_{In}) &= \sum_j P(w_{In}, t) \text{ where } t \text{ is a parse of } w_{In} \\ &= \sum_j P(t) \end{aligned}$$

# Parsing SCFG <sub>5</sub>

---

- **Positional invariance:**
  - The probability of a subtree is independent of its position in the derivation tree
- **Context-free**
- **Independence from ancestors:**

## Parameter estimation

- Supervised learning
  - From a treebank
    - $\{\psi_1, \dots, \psi_N\}$
- Non supervised learning
  - Inside/Outside (EM)
  - Similar to Baum-Welch in HMM

# Treebank grammars <sup>1</sup>

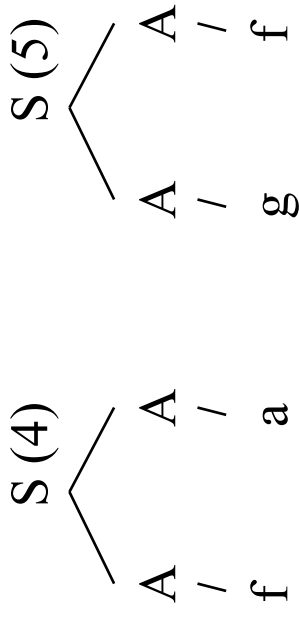
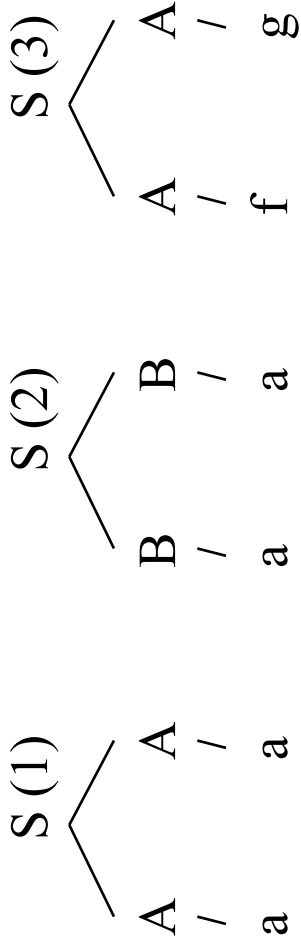
---

- Grammars directly derived from a treebank
  - Charniak, 1996
    - Using PTB
      - 47,000 sentences
    - Navigating PTB where each local subtree provides the left hand and right hand side of a rule
    - Precision and recall around 80%
    - around 17,500 rules
- Sekine, 1997, Sekine & Grishman, 1995
- Treebank grammars compactation
  - Continuous growth of the grammar size
  - Krotov et al, 1999, Krotov, 1998, Gaizauskas, 1995

# Treebank grammars <sub>2</sub>

---

- Consider a treebank containing the following trees:





## Treebank grammars <sup>3</sup>

---

- Suppose that (1) occurs 40 times, (2) occurs 10 times, (3) occurs 5 times, (4) occurs 5 times, and (5) occurs once.
- We want to induce a SCFG reflexing this treebank.
- Parameter estimation

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

# Treebank grammars <sup>4</sup>

---

- Rules

$$S \rightarrow A A \quad : \quad 40 + 5 + 5 + 1 = 51$$

$$S \rightarrow B B \quad : \quad 10$$

$$A \rightarrow a \quad : \quad 40 + 40 + 5 = 85$$

$$A \rightarrow f \quad : \quad 5 + 5 + 1 = 11$$

$$A \rightarrow g \quad : \quad 5 + 1 = 6$$

$$B \rightarrow a \quad : \quad 10$$

# Trebank grammars <sup>5</sup>

---

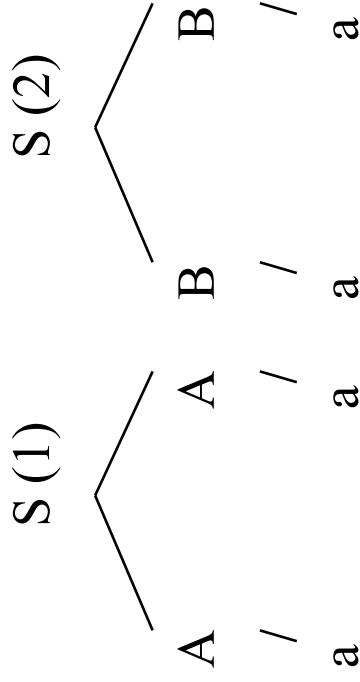
- Parameters maximizing global likelihood of the corpus:

G	Frequency	Total	Probability
$S \rightarrow AA$	51	61	0.836
$S \rightarrow BB$	10	61	0.164
$A \rightarrow a$	85	102	0.833
$A \rightarrow f$	11	102	0.108
$A \rightarrow g$	6	102	0.059
$B \rightarrow a$	10	10	1.0

## Trebank grammars <sup>6</sup>

---

- Given this parametrization, what is the most likely parse for "a a"?



- $P(1) = P(S \rightarrow A A) * P(A \rightarrow a) * P(A \rightarrow a)$   
 $= 0.836 * 0.833 * 0.833 = 0.580$
- $P(2) = P(S \rightarrow B B) * P(B \rightarrow a) * P(B \rightarrow a)$   
 $= 0.164 * 1.0 * 1.0 = 0.164$

## Treebank grammars 7

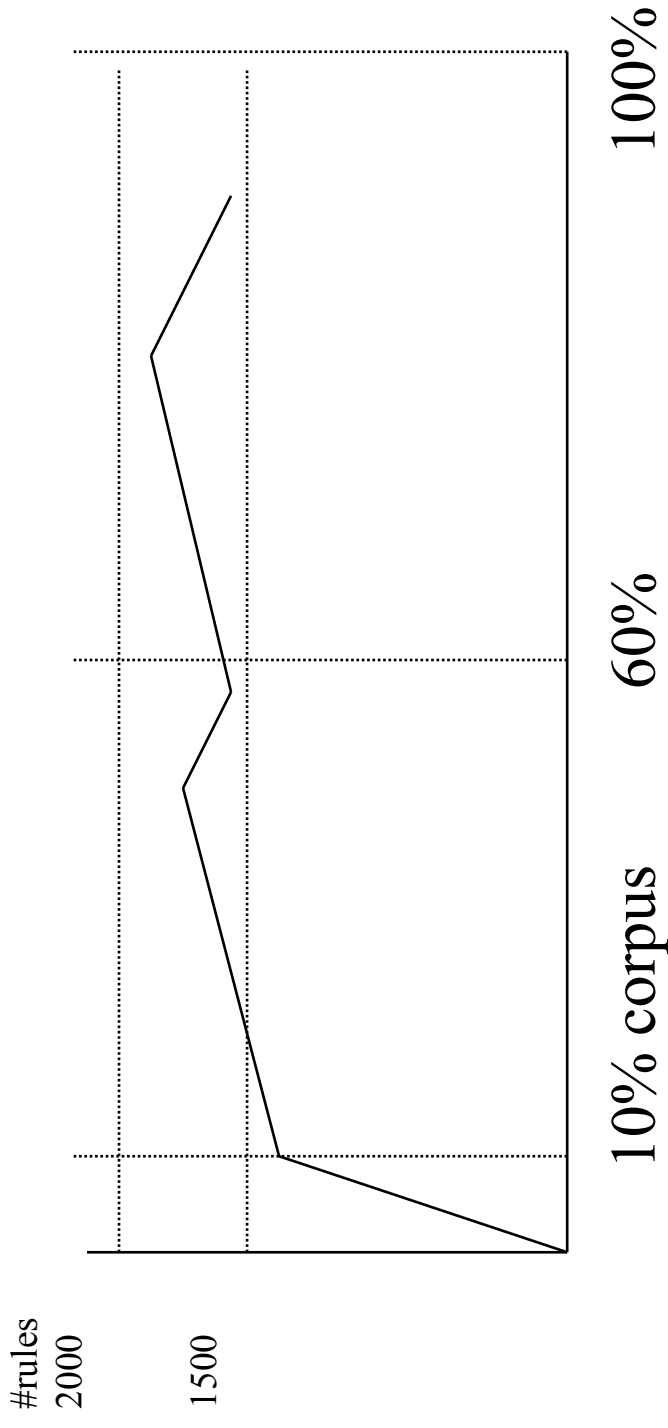
---

- Treebank Grammars compactation
- Partial bracketting
  - NP → DT NN CC DT NN
  - NP → NP CC NP
  - NP → DT NN
- Redundance removing (some rules can be generated from others)

# Treebank grammars <sup>8</sup>

---

- Applying compactation
- 17,529  $\rightarrow$  1,667 rules



# Trebank grammars 9

---

- Removing non Linguistically valid rules
  - Assign probabilities (MLE) to the initial rules
  - Remove a rule except in the the probability of the structure built from its application is greater than the probability of building applying simpler rules.
- thresholding
  - Removing rules occurring  $< n$  times

	Full	Simply thresholded	Fully compacted	Linguistically Compacted Grammar 1	Linguistically Compacted Grammar 2
Recall	70.55	70.78	30.93	71.55	70.76
precision	77.89	77.66	19.18	72.19	77.21
Grammar size	15,421	7,278	1,122	4,820	6,417

# Maximum Entropy Models (ME) <sup>1</sup>

---

- Maximum Entropy: alternative estimation technique.
- Able to deal with different kinds of evidence
- ME principle:
  - Do not assume anything about non-observed events.
  - Find the most uniform (maximum entropy, less informed) probability distribution that matches the observations.

- Example:

$p(a, b)$	0	1	$p(a, b)$	
x	?	?	x	x
y	?	?	y	y
total	0.6	1.0	total	total
			0.6	1.0

*Observations*

*One possible  $p(a, b)$*

*Max. Entropy  $p(a, b)$*



## Maximum Entropy Models (ME) <sup>2</sup>

---

- Observed facts are constraints for the desired model  $p$ .
  - Constraints take the form of feature functions:
- $$f_i : \epsilon \rightarrow \{0, 1\}$$
- The desired model must satisfy the constraints:

$$E_p(f_i) = E_{\tilde{p}}(f_i) \quad \forall i$$

where:

$$E_p(f_i) = \sum_{x \in \epsilon} p(x) f_i(x) \quad \text{expectation of model } p.$$

$$E_{\tilde{p}}(f_i) = \sum_{x \in \epsilon} \tilde{p}(x) f_i(x) \quad \text{observed expectation.}$$

# Maximum Entropy Models (ME) <sup>3</sup>

---

- Example:

$$\epsilon = \{x, y\} \times \{0, 1\}$$

$p(a, b)$	0	1
x	?	?
y	?	?
total	0.6	1.0

- Observed fact:  $p(x, 0) + p(y, 0) = 0.6$
- Encoded as a constraint:  $E_p(f_1) = 0.6$

where:

- $f_1(a, b) = \begin{cases} 1 & \text{if } b = 0 \\ 0 & \text{otherwise} \end{cases}$
- $E_p(f_1) = \sum_{(a,b) \in \{x,y\} \times \{0,1\}} p(a, b) f_1(a, b)$

# Maximum Entropy Models (ME) <sup>4</sup>

---

Ratnaparkhi ME parser

- ME based
- 4 procedures
  - Tag
  - Chunk
  - Build
  - Check
- where
- $P_X(a | b)$ 
  - X is one of the procedures
  - a is a valid action for this procedure
  - b is the context

## Maximum Entropy Models (ME) <sup>5</sup>

---

- A set of *templates* is associated to each of the procedures
- *templates* incorporates the type of features relevant for parsing
  - *Headwords* of constituents
  - Combination of *headwords*
  - Generalizations (POS, syntactic categories)
  - Limited forms of *look-ahead*

## Maximum Entropy Models (ME) <sup>6</sup>

---

- Learning is very simple:
  - Simple counting rejecting features occurring  $< 5$  times.
  - Using 40.000 sentences from *PTB* the model includes 120.000 features of type *TAG*, 230.000 of type *CHUNK*, 180.000 of type *CHECK* and 530.000 of type *BUILD* (most of them lexicalized).
- Good results:
  - recall (85.3%) precision (87.5%).

# Lexicalized parsers <sub>1</sub>

---

- Two types of statistics
- Relation between the head of a phrase and the rule applied to expand it:
  - $p(r|h)$
  - $r = \text{rule}, h = \text{head}$
- Relation between the head of a phrase and another sub-component:
  - $p(h|m,t)$
  - $h = \text{head of the sub-component}, m = \text{head of the component}, t = \text{type of the sub-component}$
- probability of a lexicalized parse tree:
  - $p(s,\pi) = p(h(c)|m(c),t) \cdot p(r(c)|h(c))$

## Lexicalized parsers <sup>2</sup>

---

- Other useful statistics
- The type of the parent conditions the probability of the rule (Charniak, 1997)
- Immediate left context (Collins, 1996)
- Classification of NP in a VP as optional or mandatory (Collins, 1997)
- ...

## Lexicalized parsers <sup>3</sup>

---

Collins parser

- parser based in SCFG.

$$\psi_{\text{OPT}} = \underset{\psi}{\operatorname{argmax}} P(\psi | s) = \underset{\psi}{\operatorname{argmax}} P(\psi, s)$$

- where  $\psi_{\text{OPT}}$  is the most probable parser and  $\psi$  represents whatever parse tree of sentence  $s$ . If  $\psi$  is obtained applying  $n$  rules of type  $A \rightarrow \alpha$  for producing  $s$ , we get:

$$P(\psi, s) = \prod_{i=1}^n P((A \rightarrow \alpha)_i) = \prod_{i=1}^n P(\alpha_i | A_i)$$



## Lexicalized parsers <sup>4</sup>

---

- The rules of the grammar are of the form that for each rule the left side has attached a word (*head*) taken from one of the components of the right side.
- $P$  parent component,  $H$  head,  $L_i$  and  $R_j$  components to the left (left modifiers) and right (right modifiers) of the head,  $h$ , *head*,  $l_i$  and  $r_j$  words associated to these components:
  - $P(h) \rightarrow L_n(l_n) \dots L_1(l_1) H(h) R_1(r_1) \dots R_m(r_m)$

## Lexicalized parsers <sup>5</sup>

---

- The initial model, lexicalized, cannot be scaled up and Collins decompose it into 3 simpler models which can be learned from the *PTB* and can generate the right side part
- The first generates  $H$  with probability  $P(H|P, h)$  given a category  $P$  and a headword  $h$ .

## Lexicalized parsers <sup>5</sup>

---

- Then the right modifiers at distance  $i = 1, \dots, m+1$  (where  $R_{m+1}(r_{m+1}) = \text{STOP}$ ):

$$\prod_{i=1}^{m+1} P_r(R_i(r_i) | P, h, H)$$

- Now conditioned not only to  $P$  and  $h$  but also to  $H$
- Finally the left modifiers:

$$\prod_{i=1}^{n+1} P_l(L_i(l_i) | P, h, H)$$

## Lexicalized parsers <sup>6</sup>

---

- For instance, the probability associated to the following rule:
  - $S(\text{bought}) \rightarrow NP(\text{week}) NP(\text{Mark}) VP(\text{bought})$
  - would be:

$$P_h(VP | S, \text{bought}) \cdot P_1(NP(\text{Mark}) | S, \text{bought}, VP) \cdot$$

$$P_1(NP(\text{week}) | S, \text{bought}, VP) \cdot P_1(\text{STOP} | S, \text{bought}, VP) \cdot$$

$$P_h(\text{STOP} | S, \text{bought}, VP) \cdot$$

## Lexicalized parsers <sup>7</sup>

---

- The parser was learned with 40.000 sentences of PTB
- Results:
  - 87.4% recall and 88.1% precision.
- In Collins (1997) the model is improved with linguistic information on the distinction between complements and adjuncts (model 2) and constituent moving, relative sentences (modelo 3).
- With these improvements the parser reaches 88.1% of recall and 88.6% of precision.