



max planck institut
informatik

PCFG: Probabilistic Context Free Grammars

Presenter: Ba Dat Nguyen

Advisor: Dr. Martin Theobald

Max-Planck-Institut für Informatik
Saarbrücken, Germany

Outline

- ➔ • Introduction
- **Probabilistic Context Free Grammars**
 - Parsing
 - **Context Free Grammars**
 - **Probabilistic Context Free Grammars**
 - Inside-Outside Algorithm
- **Extension**
 - Distance
 - Complement/ adjunct distinction
 - Traces and Wh-movement

The World is a big ambiguity



Solution

PCFG is a good way to solve ambiguity problems in syntactic structure field.

Outline

- Introduction
- Probabilistic Context Free Grammars
 - ➔ Parsing
 - Context Free Grammars
 - Probabilistic Context Free Grammars
 - Inside-Outside Algorithm
- Extension
 - Distance
 - Complement/ adjunct distinction
 - Traces and Wh-movement

Language and Grammar

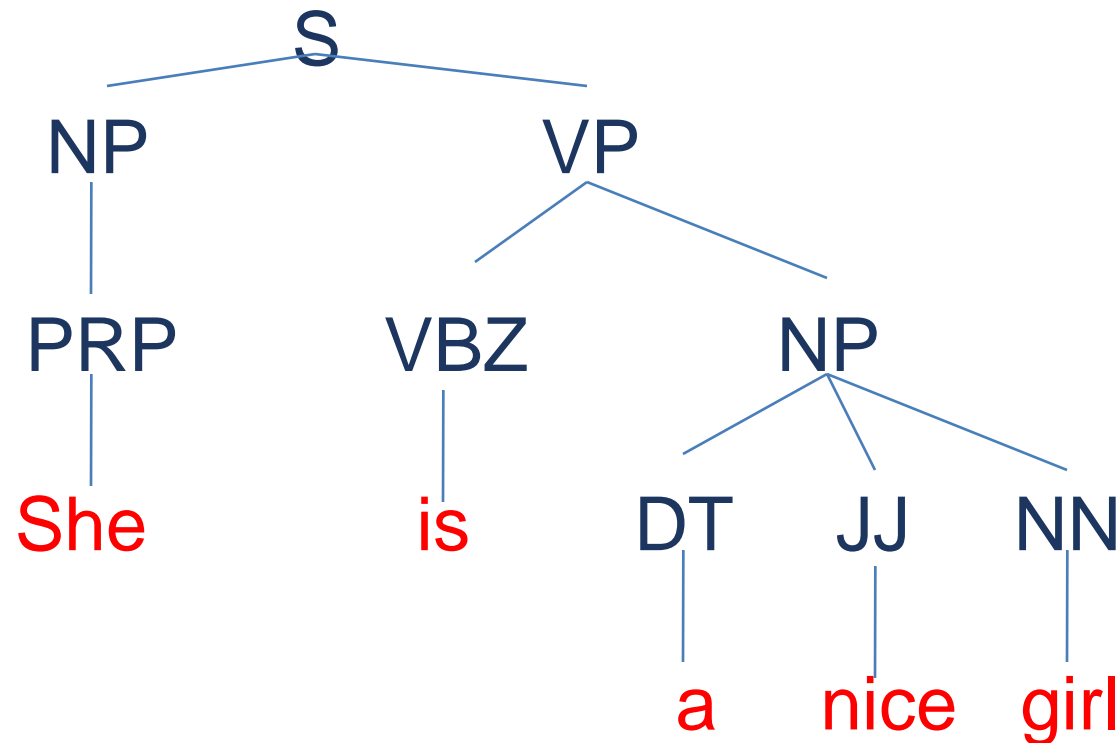
- Language
 - Structural
 - Ambiguous
- Grammar
 - Generalization of regularities in language structures
 - Morphology and syntax

Parsing

- Process working out the grammatical structure of sentences.
- Basic Parsing Algorithms
 - Parsing Strategies
 - CYK Algorithm
 - Earley Algorithm

Example of parsing

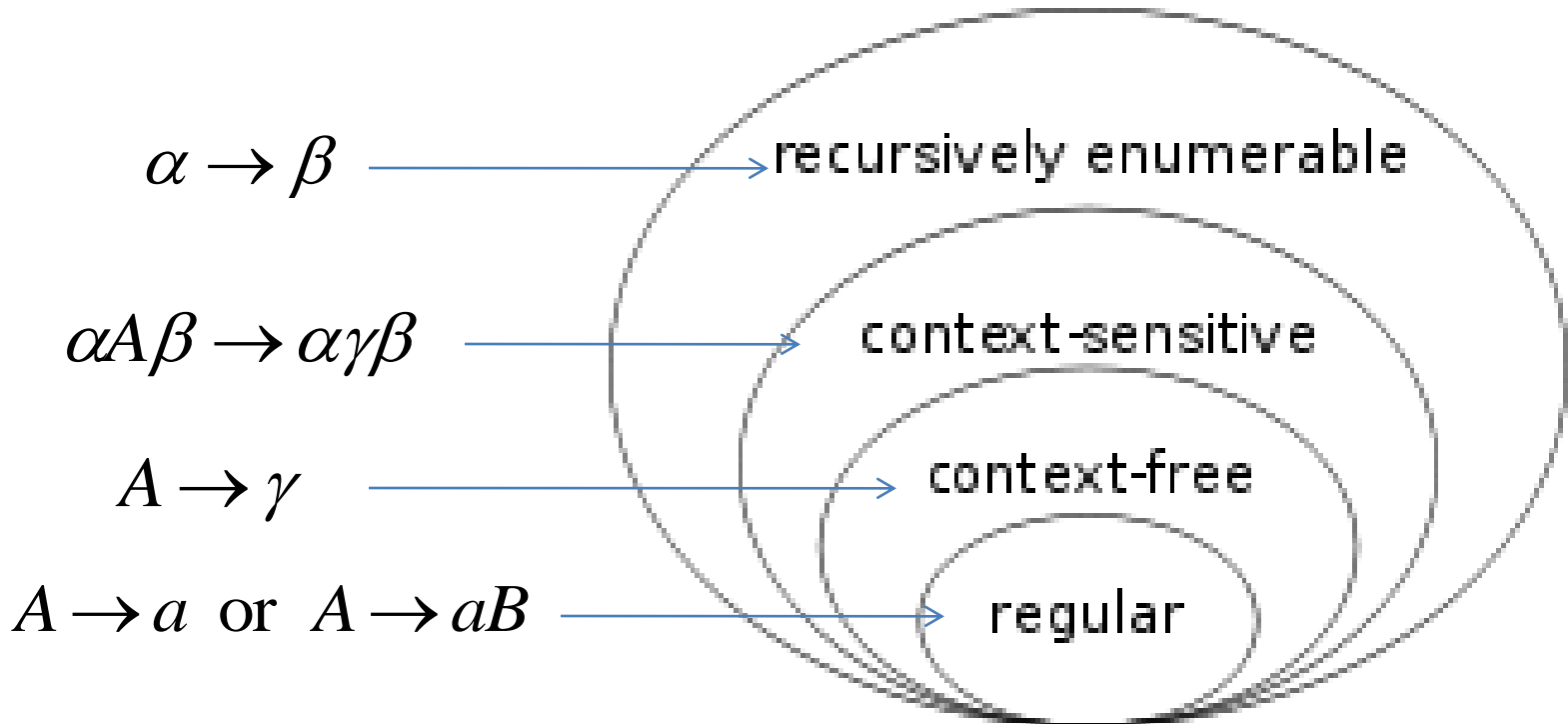
- “She is a nice girl”



Outline

- Introduction
- Probabilistic Context Free Grammars
 - Parsing
 - ➔ ▪ **Context Free Grammars**
 - **Probabilistic Context Free Grammars**
 - Inside-Outside Algorithm
- Extension
 - Distance
 - Complement/ adjunct distinction
 - Traces and Wh-movement

Chomsky hierarchy



Where :

A, B are nonterminals

a is a terminal

α, β, γ are strings of terminals and nonterminals

Context Free Grammars (CFG)

- A **C**ontext **F**ree **G**rammars consists of
 - 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 $\{N^i \rightarrow \xi^j\}$
where ξ^j is a sequence of terminals and nonterminals

Example of CFG

S \rightarrow NP VP

PP \rightarrow P NP

VP \rightarrow V NP

VP \rightarrow VP PP

P \rightarrow with

V \rightarrow saw

NP \rightarrow NP PP

NP \rightarrow astronomers

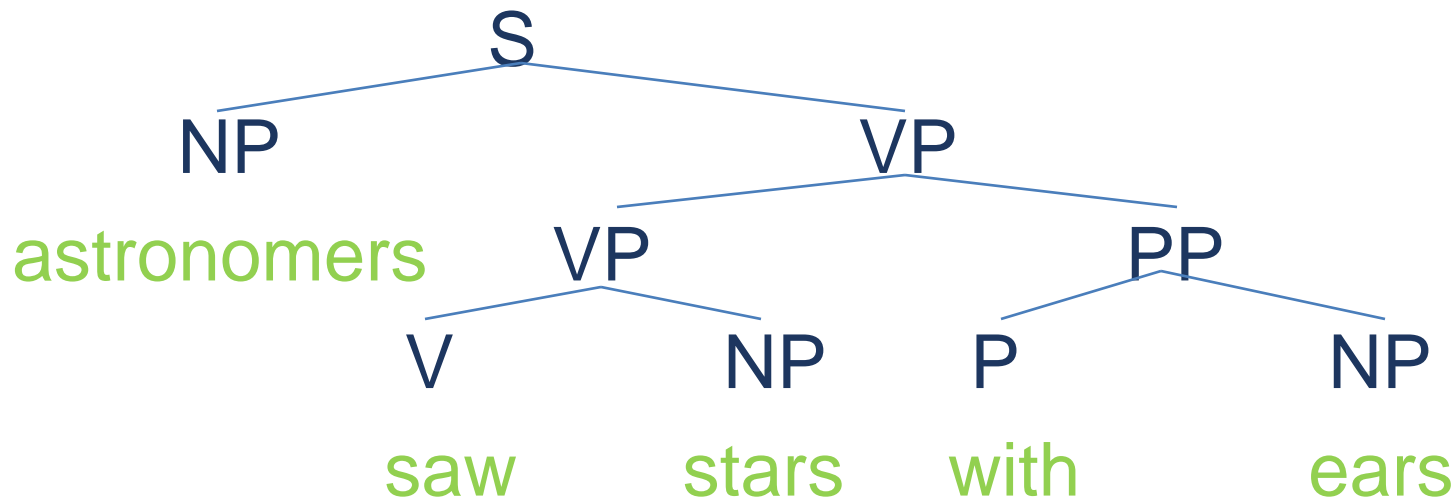
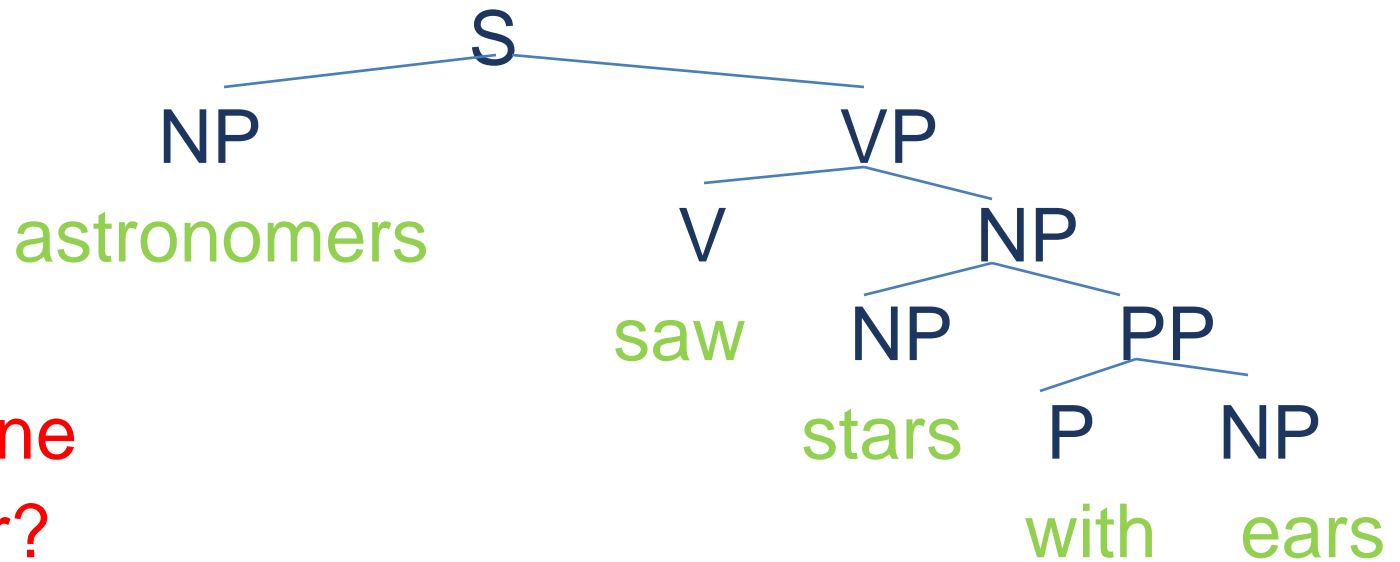
NP \rightarrow ears

NP \rightarrow saw

NP \rightarrow stars

NP \rightarrow telescopes

Ambiguous sentences



Outline

- Introduction
- Probabilistic Context Free Grammars
 - Parsing
 - Context Free Grammars
 - ➔ ▪ **Probabilistic Context Free Grammars**
 - Inside-Outside Algorithm
- Extension
 - Distance
 - Complement/ adjunct distinction
 - Traces and Wh-movement

Probabilistic CFG

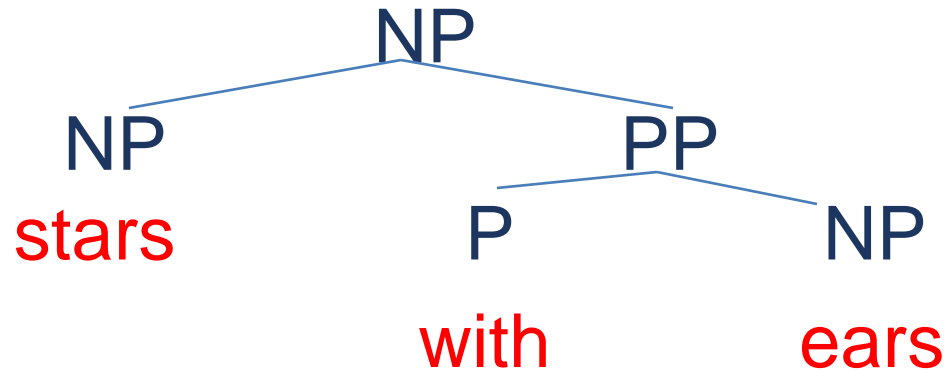
- A **P**robabilistic **C**ontext **F**ree **G**rammars (PCFG) consists of
 - A CFG
 - A corresponding set of probabilities on rules such that:

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

Example of PCFG

S -> NP VP	1.0	NP -> NP PP	0.4
PP -> P NP	1.0	NP -> astronomers	0.1
VP -> V NP	0.7	NP -> ears	0.18
VP -> VP PP	0.3	NP -> saw	0.04
P -> with	1.0	NP -> stars	0.18
V -> saw	1.0	NP -> telescopes	0.1

Probability of a tree

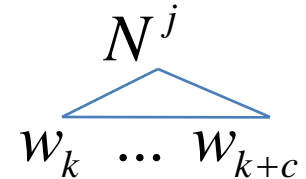


$$\begin{aligned}
 & P(NP \rightarrow NP PP, NP \rightarrow stars, PP \rightarrow P NP, P \rightarrow with, NP \rightarrow ears) \\
 &= P(S \rightarrow NP PP) \times P(NP \rightarrow stars \mid NP \rightarrow NP PP) \\
 &\quad \times P(PP \rightarrow P NP \mid NP \rightarrow NP PP, NP \rightarrow stars) \\
 &\quad \times P(P \rightarrow with \mid PP \rightarrow P NP, NP \rightarrow NP PP, NP \rightarrow stars) \\
 &\quad \times P(NP \rightarrow ears \mid P \rightarrow with, PP \rightarrow P NP, NP \rightarrow NP PP, NP \rightarrow stars)
 \end{aligned}$$

Assumptions

- Place invariance

$\forall k \ P(N_{k(k+c)}^j \rightarrow \xi)$ *is the same*



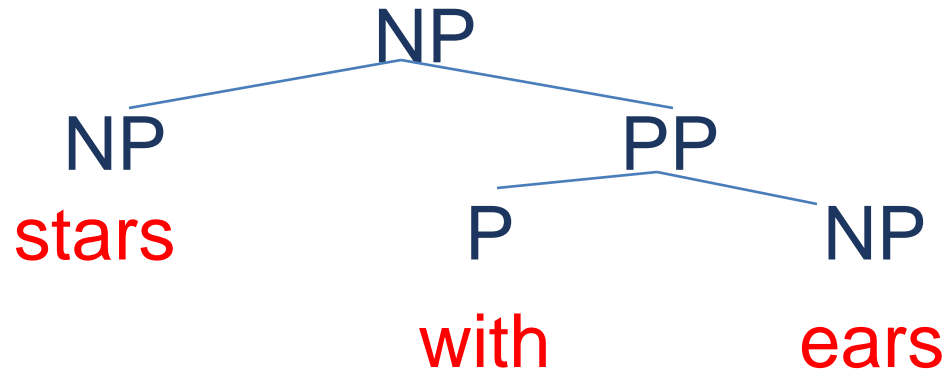
- Context-free

$P(N_{kl}^j \rightarrow \xi \mid \text{anything outside } k \text{ through } l) = P(N_{kl}^j \rightarrow \xi)$

- Ancestor-free

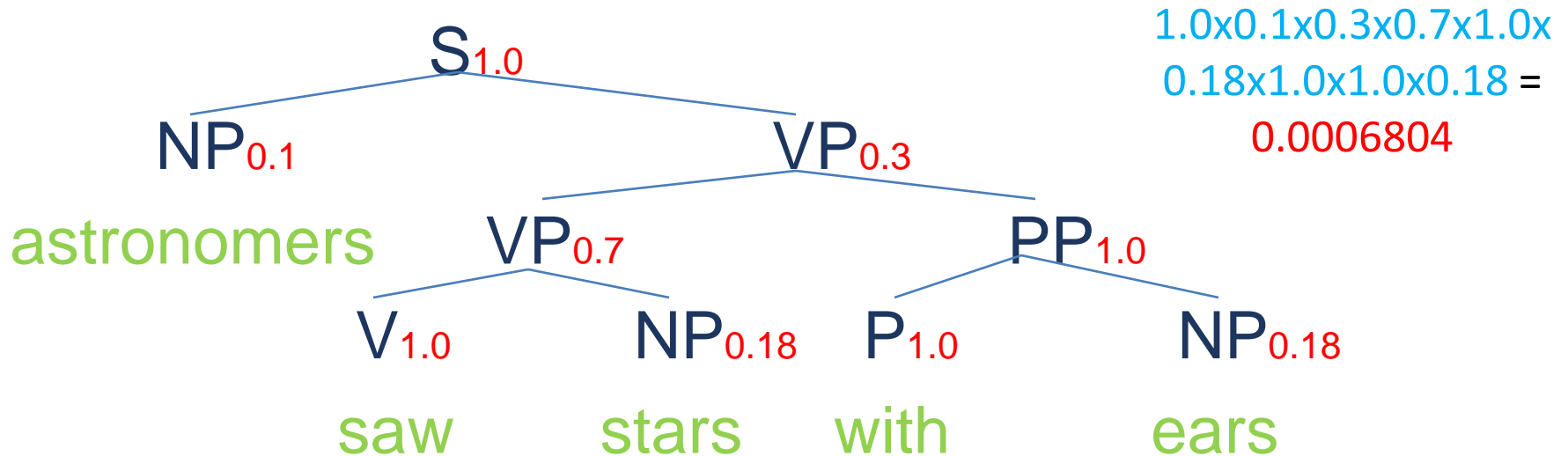
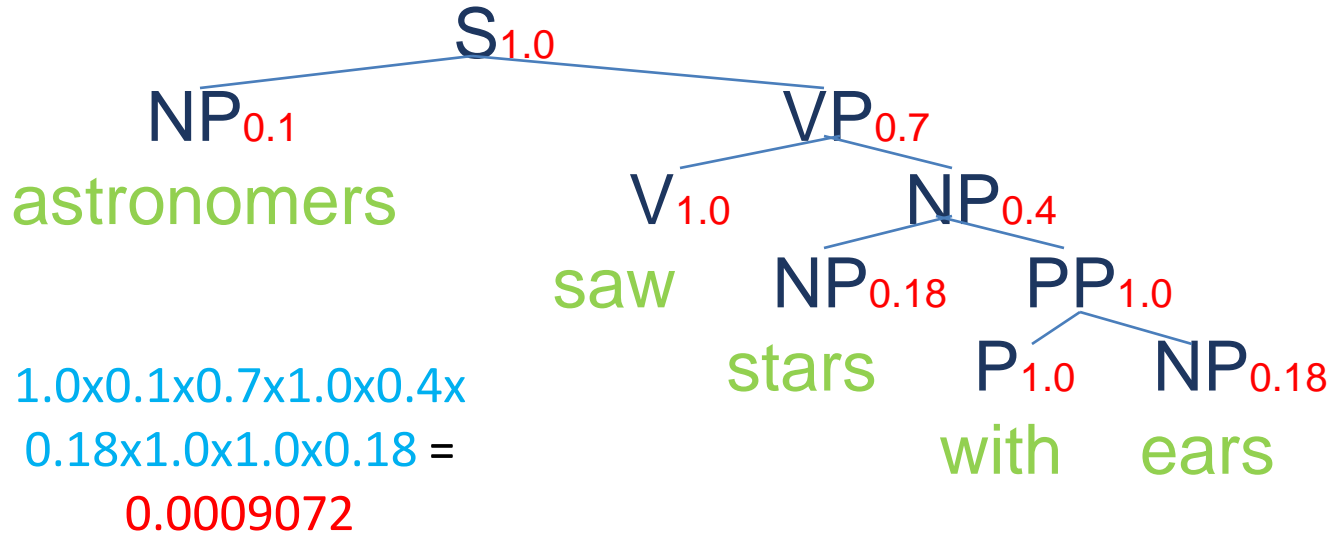
$P(N_{kl}^j \rightarrow \xi \mid \text{any ancestor nodes outside } N_{kl}^j) = P(N_{kl}^j \rightarrow \xi)$

Probability of a tree



$$\begin{aligned}
 &P(NP \rightarrow NP PP, NP \rightarrow stars, PP \rightarrow P NP, P \rightarrow with, NP \rightarrow ears) \\
 &= P(S \rightarrow NP PP) \times P(NP \rightarrow stars) \times P(PP \rightarrow P NP) \\
 &\quad \times P(P \rightarrow with) \times P(NP \rightarrow ears)
 \end{aligned}$$

Ambiguity



Outline

- Introduction
- Probabilistic Context Free Grammars
 - Parsing
 - Context Free Grammars
 - Probabilistic Context Free Grammars
 - ➔ ▪ Inside-Outside Algorithm
- Extension
 - Distance
 - Complement/ adjunct distinction
 - Traces and Wh-movement

Probability of a rule

- Given a training set of annotated sentences

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

C(.) - number of times that a particular rule is used.

Probability of a rule

How to calculate if
there is no
annotated data!



Maximum Likelihood Estimation

- Maximum Likelihood Estimation

$$\arg \max_{\mu} P(O_{training} | \mu)$$

μ = parameters of current grammar set

- No known analytic method to choose μ to maximize $P(O | \mu)$
- Locally maximize $P(O | \mu)$ by an iterative hill-climbing – special case of **Expectation Maximization** method.
- Inside-Outside algorithm is a form of EM using the inside-outside probabilities estimated from training set.

Training a PCFG

- We are given
 - A set of training sentences
 - A set of terminals
 - A set of nonterminals
- Initial probabilities are estimated by rules (perhaps by randomly chosen)
- Using inside-outside algorithm to train

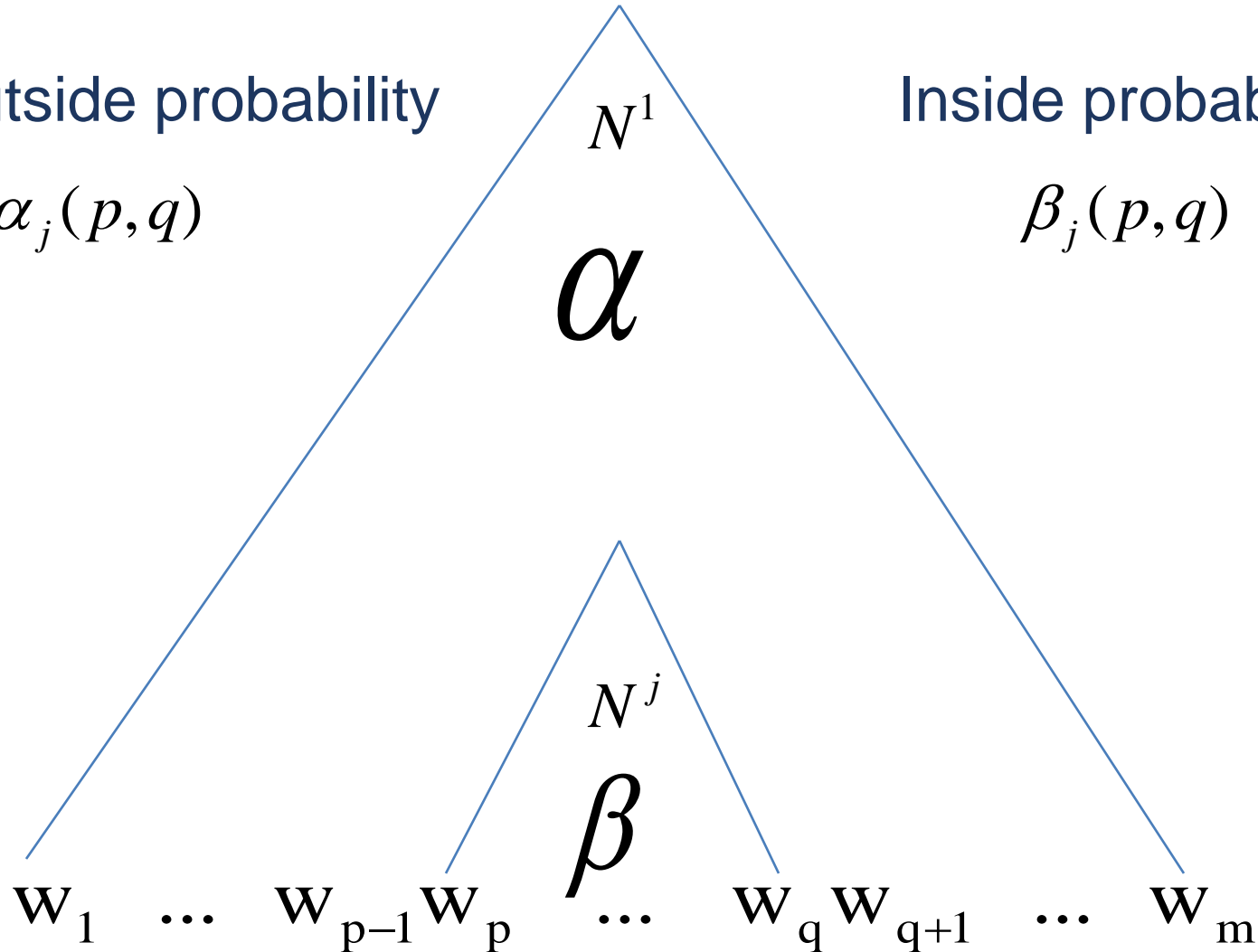
Inside-Outside probabilities

- Outside probability

$$\alpha_j(p, q)$$

Inside probability

$$\beta_j(p, q)$$



Inside probabilities

- Inside probability $\beta_j(p, q)$ is the probability of sequence $W_p \dots W_q$ being generated with a tree rooted by node N^j

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

- Calculation can be carried out bottom-up

$$\beta_j(k, k) = P(N^j \rightarrow w_k)$$

$$\beta_j(p, q) = \sum_{r,s} \sum_{d=q}^{q-1} P(N^j \rightarrow N^r N^s) \beta_r(p, d) \beta_s(d+1, q)$$

Outside probabilities

- Outside probability $\alpha_j(p, q)$ is the total probability of beginning with the start symbol and generating all the words outside N_{pq}^j

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

$$N_{pq}^j = N^j \xrightarrow{*} w_p \dots w_q$$

$$\alpha_1(1, m) = 1 \text{ and } \alpha_j(1, m) = 0 \text{ for } j \neq 1$$

$$\begin{aligned} \alpha_j(p, q) = & \sum_{f, g} \sum_{e=q+1}^m \alpha_f(p, e) P(N^f \rightarrow N^j N^g) \beta_g(q+1, e) \\ & + \sum_{f, g} \sum_{e=1}^{p-1} \alpha_f(e, q) P(N^f \rightarrow N^g N^j) \beta_g(e, p-1) \end{aligned}$$

Inside-Outside Algorithm

We have:

$$\alpha_j(p, q) \beta_j(p, q) = P(N^1 \xrightarrow{*} w_{1,m}, N^j \xrightarrow{*} w_{p,q})$$

$$\text{Call } P(N^1 \xrightarrow{*} w_{1,m}) \pi$$

$$P(N^j \xrightarrow{*} w_{p,q} \mid N^1 \xrightarrow{*} w_{1,m}) = \frac{\alpha_j(p, q) \beta_j(p, q)}{\pi}$$

Inside-Outside Algorithm

$$E(N^j \text{ is used}) = \frac{\sum_{p=1}^m \sum_{q=p}^m \alpha_j(p, q) \beta_j(p, q)}{\pi}$$

$$E(N^j \rightarrow N^r N^s, N^j \text{ used})$$

$$= \frac{\sum_{p=1}^{m-1} \sum_{q=p+1}^m \sum_{d=p}^{q-1} \alpha_j(p, q) P(N^j \rightarrow N^r N^s) \beta_r(p, d) \beta_s(d+1, q)}{\pi}$$

Inside-Outside Algorithm

Therefore:

$$\begin{aligned}
 & P(N^j \rightarrow N^r N^s) \\
 &= \frac{\sum_{p=1}^{m-1} \sum_{q=p+1}^m \sum_{d=p}^{q-1} \alpha_j(p, q) P(N^j \rightarrow N^r N^s) \beta_r(p, d) \beta_s(d+1, q)}{\sum_{p=1}^m \sum_{q=1}^m \alpha_j(p, q) \beta_j(p, q)} \\
 & P(N^j \rightarrow w^k) = \frac{\sum_{h=1}^m \alpha_j(h, h) P(w_h = w^k) \beta_j(h, h)}{\sum_{p=1}^m \sum_{q=1}^m \alpha_j(p, q) \beta_j(p, q)}
 \end{aligned}$$

Inside-Outside Algorithm

For each sentence W^i in the corpus

$$f_i(p, q, j, r, s) =$$

$$\frac{\sum_{d=p}^{q-1} \alpha_j(p, q) P(N^j \rightarrow N^r N^s) \beta_r(p, d) \beta_s(d+1, q)}{P(N^1 \xrightarrow{*} W_i)}$$

$$g_i(h, j, k) = \frac{\alpha_j(h, h) P(w_h = w^k) \beta_j(h, h)}{P(N^1 \xrightarrow{*} W_i)}$$

$$h_i(p, q, j) = \frac{\alpha_j(p, q) \beta_j(p, q)}{P(N^1 \xrightarrow{*} W_i)}$$

Inside-Outside Algorithm

We have

$$P(N^j \rightarrow N^r N^s) = \frac{\sum_{i=1}^l \sum_{p=1}^{m_i-1} \sum_{q=p+1}^{m_i} f_i(p, q, j, r, s)}{\sum_{i=1}^l \sum_{p=1}^{m_i} \sum_{q=p}^{m_i} h_i(p, q, j)}$$

$$P(N^j \rightarrow w^k) = \frac{\sum_{i=1}^l \sum_{h=1}^{m_i} g_i(h, j, k)}{\sum_{i=1}^l \sum_{p=1}^{m_i} \sum_{q=p}^{m_i} h_i(p, q, j)}$$

Discussion

- Inside-Outside algorithm is quite slow $O(m^3 n^3)$ for each sentence
 - m is the length of the sentence
 - n is the number of nonterminals
- The algorithm is very sensitive to the initialization of the parameters.
- In practice, a PCFG is a worse language model for English than an n -gram model ($n > 1$).

Outline

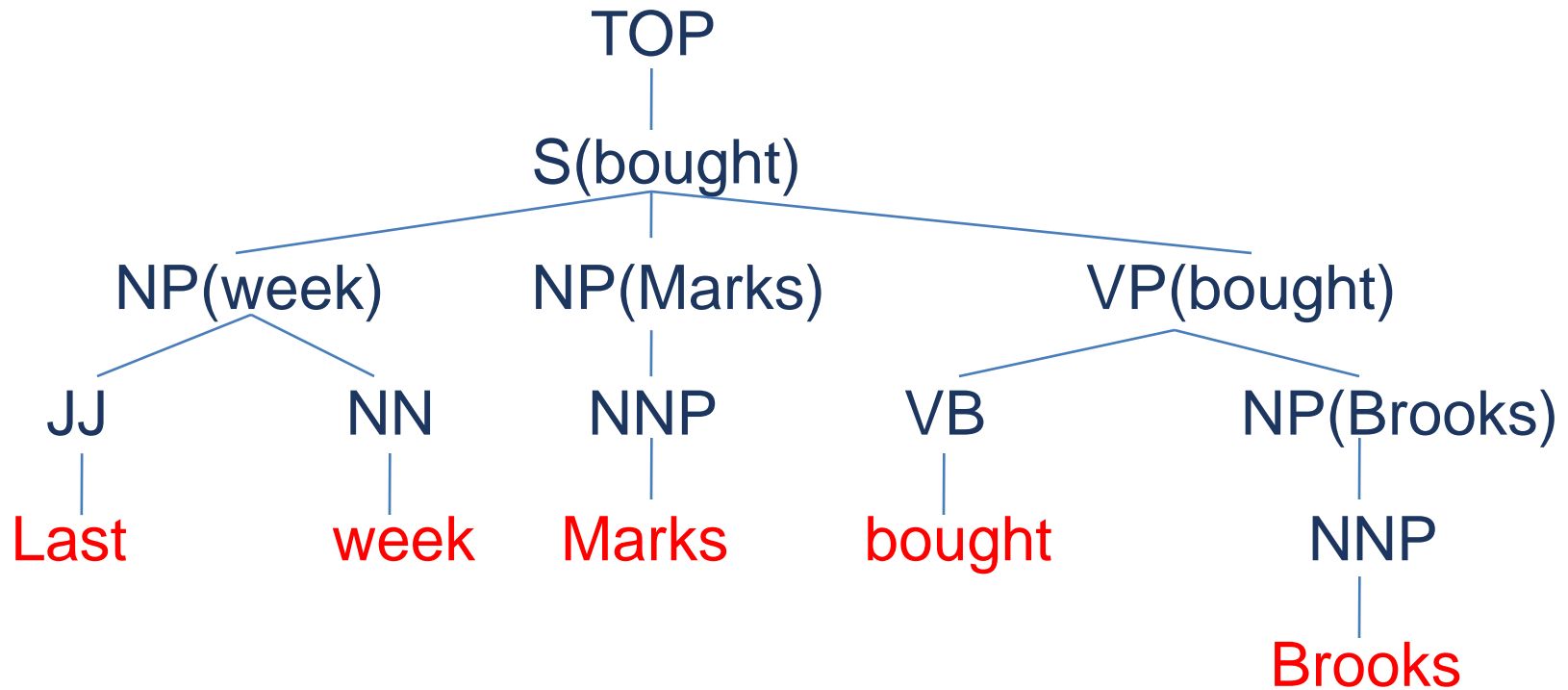
- Introduction
- Probabilistic Context Free Grammars
 - Parsing
 - Context Free Grammars
 - Probabilistic Context Free Grammars
 - Inside-Outside Algorithm
- **Extension**
 - ➔ ▪ **Distance**
 - Complement/ adjunct distinction
 - Traces and Wh-movement

More features

TOP	->	S(bought)		
S(bought)	->	NP(week)	NP(Marks)	VP(bought)
NP(week)	->	JJ>Last)	NN(week)	
NP(Marks)	->	NNP(Marks)		
VP(bought)	->	VB(bought)	NP(Brooks)	
NP(Brooks)	->	NNP(Brooks)		

- Adding some lexical features: words and POS inside non-terminals.
- Using the head-child of the phrase, which inherits the head-word from its parent.

Example



Using Distance

- Using some probabilities
 - Head constituent label of the phrase

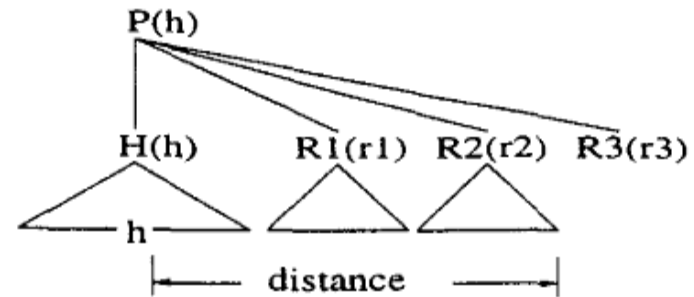
$$P_H(H | P, h)$$

- Modifiers to the right of the head

$$\prod_{i=1..m+1} P_R(R_i(r_i) | P, h, H, R_1(r_1) \dots R_{i-1}(r_{i-1}))$$

- Modifiers to the left of the head

$$\prod_{i=1..n+1} P_L(L_i(l_i) | P, h, H, L_1(l_1) \dots L_{i-1}(l_{i-1}))$$



Using Distance

For example :

$$\begin{aligned}
 &P(S(\text{bought}) \rightarrow NP(\text{week})NP(\text{Marks})VP(\text{bought})) \\
 &= P_h(VP \mid S, \text{bought}) \times P_l(NP(\text{Marks}) \mid S, VP, \text{bought}) \\
 &\times P_l(NP(\text{week}) \mid S, VP, \text{bought}, NP, \text{Marks})
 \end{aligned}$$

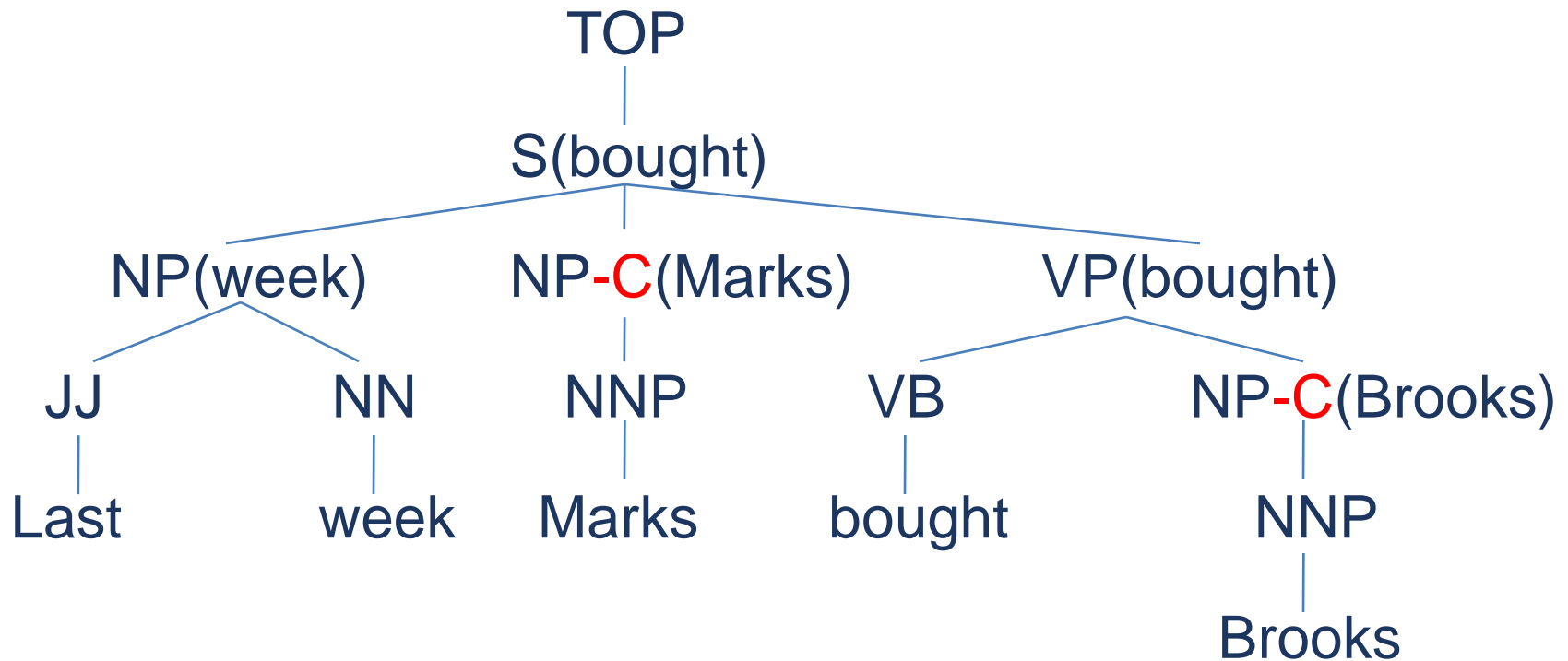
With distance 0

$$\begin{aligned}
 &P(S(\text{bought}) \rightarrow NP(\text{week})NP(\text{Marks})VP(\text{bought})) \\
 &= P_h(VP \mid S, \text{bought}) \times P_l(NP(\text{Marks}) \mid S, VP, \text{bought}) \\
 &\times P_l(NP(\text{week}) \mid S, VP, \text{bought})
 \end{aligned}$$

Outline

- Introduction
- Probabilistic Context Free Grammars
 - Parsing
 - Context Free Grammars
 - Probabilistic Context Free Grammars
 - Inside-Outside Algorithm
- **Extension**
 - Distance
 - **Complement/ adjunct distinction**
 - Traces and Wh-movement

Adding the complement / adjunct distinction



It would be useful to identify “Marks” as a subject and “Last week” as an adjunct!

Adding the complement / adjunct distinction

- Adding “-C” suffix to all non-terminals in training data which satisfy:
 - The nonterminal must be: an NP, SBAR or S whose parent is an S; an NP, SBAR, S, or VP whose parent is a VP; or S whose parent is an SBAR
 - The non-terminal must not have one of the following semantic tags: ADV, VOC, BNF, DIR, EXT, LOC, MNR, TMP, CLR or PRP.

Adding the complement / adjunct distinction

- Using some probabilities
 - Head constituent label of the phrase $P_H(H | P, h)$
 - Left and right subcat frames $P_{lc}(LC | P, H, h)$ and $P_{rc}(RC | P, H, h)$
 - Modifiers to the right of the head $P_R(R_i, r_i | H, P, h, \text{distance}(i-1), RC)$
 - Modifiers to the left of the head $P_L(L_i, l_i | H, P, h, \text{distance}(i-1), LC)$

Adding the complement / adjunct distinction

- For example

$$\begin{aligned}
 &P(S(\text{bought}) \rightarrow NP(\text{week}) NP - C(\text{Marks}) VP(\text{bought})) \\
 &= P_h(VP | S, \text{bought}) \\
 &\times P_{lc}(\{NP - C\} | S, VP, \text{bought}) \\
 &\times P_l(\{NP - C(\text{Marks})\} | S, VP, \text{bought}, \{NP - C\}) \\
 &\times P_l(NP(\text{week}) | S, VP, \text{bought})
 \end{aligned}$$

Outline

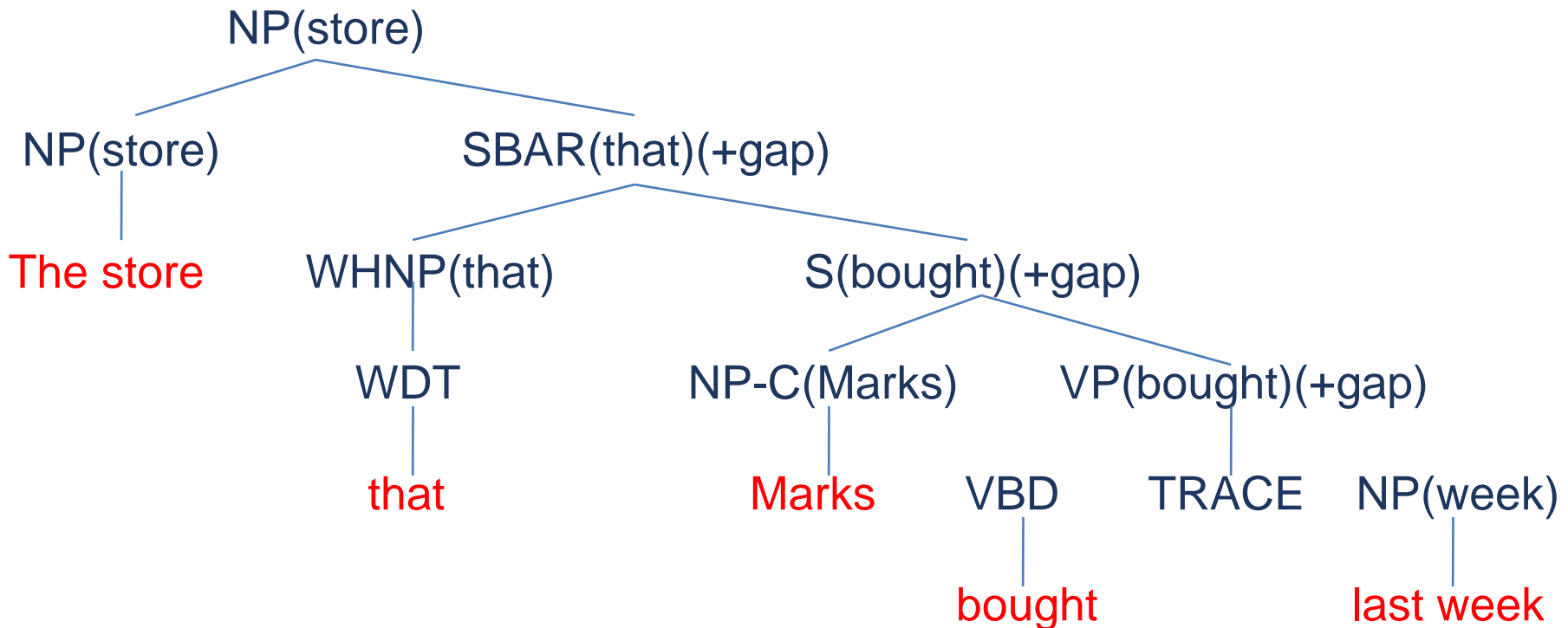
- Introduction
- Probabilistic Context Free Grammars
 - Parsing
 - Context Free Grammars
 - Probabilistic Context Free Grammars
 - Inside-Outside Algorithm
- **Extension**
 - Distance
 - Complement/ adjunct distinction
 - ➔ ▪ **Traces and Wh-movement**

Traces and Wh-movement

- Adding a “*gap*” feature to each non-terminal in the tree and propagating gaps through the tree until they are finally discharged as a trace complement.
- For example

(1)	NP	->	NP	SBAR(+gap)
(2)	SBAR(+gap)	->	WHNP	S-C(+gap)
(3)	S(+gap)	->	NP-C	VP(+gap)
(4)	VP(+gap)	->	VB	TRACE NP

Traces and Wh-movement



$$P_l(VP(bought)(+gap) \rightarrow VB(bought) TRACE NP(week))$$

$$= P_h(VB | VP, bought) \times P_G(Right | VP, bought, VB)$$

$$\times P_{RC}(\{NP - C\} | VP, bought, VB)$$

$$\times P_R(TRACE | VP, bought, VB, \{NP - C, +gap\})$$

$$\times P_R(NP(week) | VP, bought, VB)$$

Experiment

Models are trained on sections 02 - 21 of the Wall Street Journal portion of the Penn Treebank and tested on section 23.

Labelled Precision =

$$\frac{\text{number of correct constituents in proposed parse}}{\text{number of constituents in proposed parse}}$$

Labelled recall =

$$\frac{\text{number of correct constituents in proposed parse}}{\text{number of constituents in treebank parse}}$$

Crossing Brackets = number of constituents which violate constituent boundaries with a constituent in the treebank parse

MODEL	≤ 40 Words (2245 sentences)					≤ 100 Words (2416 sentences)				
	LR	LP	CBs	0 CBs	≤ 2 CBs	LR	LP	CBs	0 CBs	≤ 2 CBs
(Magerman 95)	84.6%	84.9%	1.26	56.6%	81.4%	84.0%	84.3%	1.46	54.0%	78.8%
(Collins 96)	85.8%	86.3%	1.14	59.9%	83.6%	85.3%	85.7%	1.32	57.2%	80.8%
Model 1	87.4%	88.1%	0.96	65.7%	86.3%	86.8%	87.6%	1.11	63.1%	84.1%
Model 2	88.1%	88.6%	0.91	66.5%	86.9%	87.5%	88.1%	1.07	63.9%	84.6%
Model 3	88.1%	88.6%	0.91	66.4%	86.9%	87.5%	88.1%	1.07	63.9%	84.6%

References

1. Foundations of Statistical Natural Language Processing
2. Stanford parse
<http://nlp.stanford.edu:8080/parser/>
3. Three Generative, Lexicalised Models for Statistical Parsing, ACL 97

Thanks!