# Probabilistic Context Free Grammars

## Chapter 11 from the book
### *Foundations of Statistical Natural Language Processing*
by Christopher D. Manning and Hinrich Schütze.

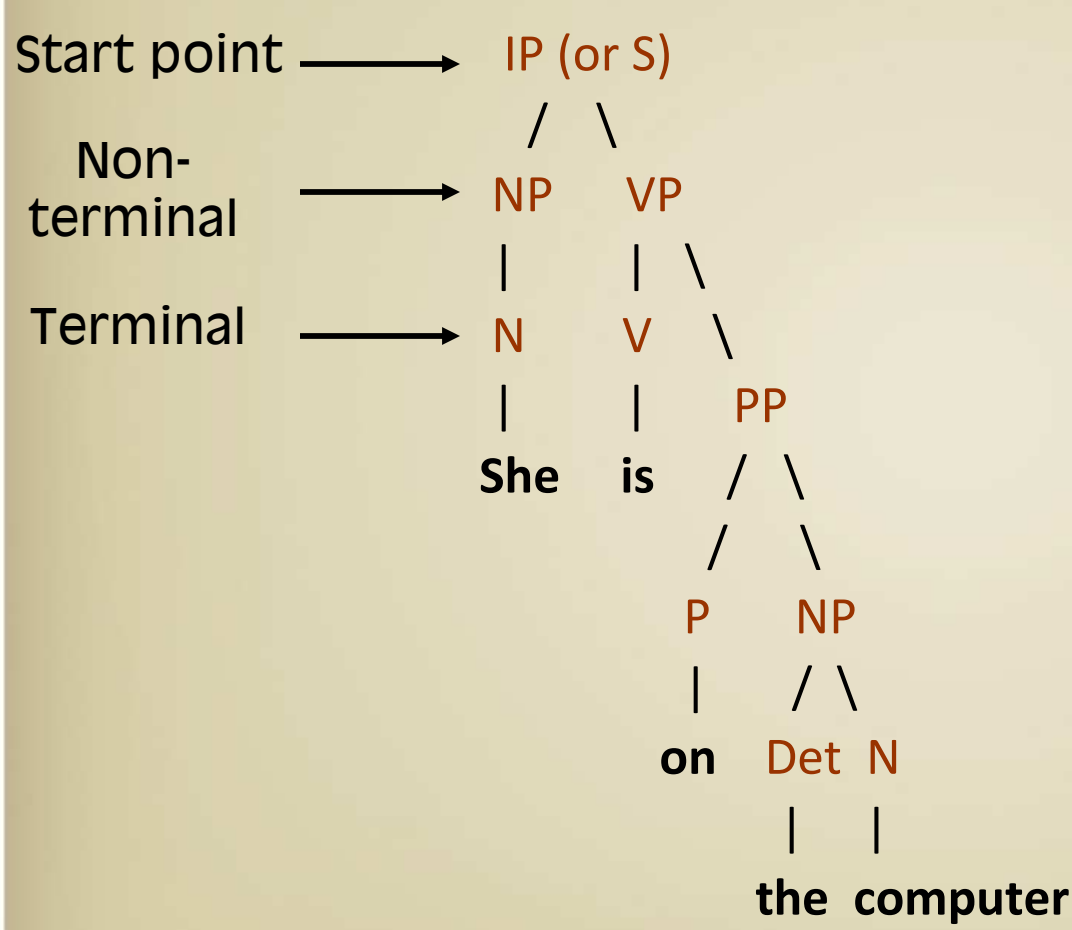## Presentation by Federico Cirett

# Outline

- Introduction

- Probabilistic Context Free Grammars (PCFG)

- Questions for PCFGs
  - Probability of a sentence
  - Most likely parse for a sentence
  - Choose a rule to maximize the prob. of a sentence

- Training a PCFG
  - Inside-Outside algorithm

# Introduction

- The quest for finding structure in language
  - Linguistics
  - Noam Chomsky 1950's - 1960's CFGs
  - Booth and Thomson 1969-1973 & others
- Uses of PCFGs
  - Speech recognition
  - Optical character recognition
  - Word grammar checker
  - Automatic translation
  - DNA sequencing

# Quick Example of parse tree (CFG)

```
Start point  ────────▶   IP (or S)
                          /  \
  Non-                  NP    VP
  terminal  ─────────▶  |     |  \
                        N     V   \
  Terminal  ─────────▶  |     |    PP
                        |     |    / \
                       She    is  /   \
                                 /     \
                                P       NP
                                |      /  \
                               on    Det   N
                                      |    |
                                     the  computer
```

- S = Start point
- IP = Inflectional phrase (sentence)
- NP = Noun phrase
- N = Noun
- VP = Verb phrase
- PP = Prepositional phrase
- P = Preposition
- Det = Determiner

*Probabilistic Context Free Grammars*

4

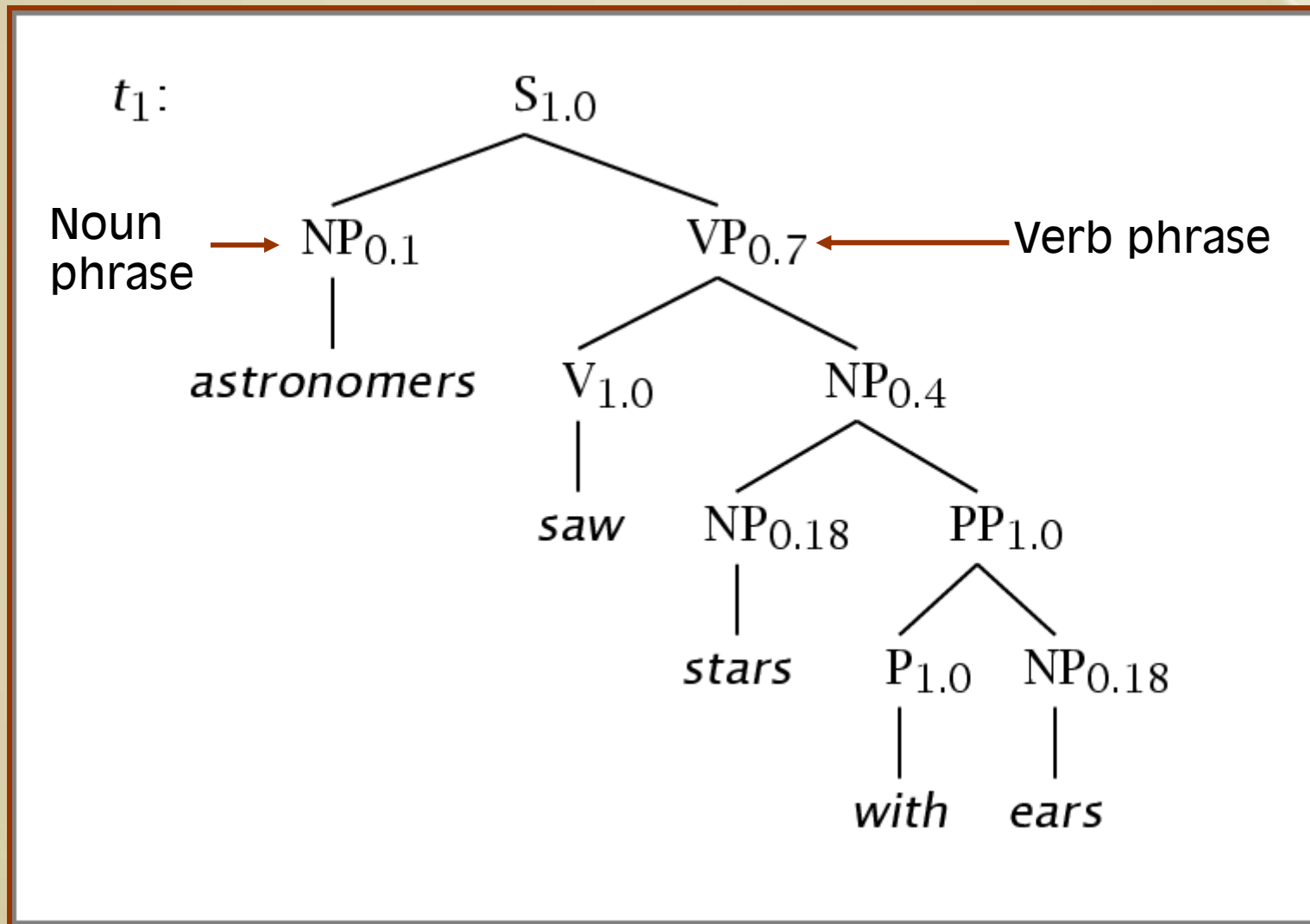# Probabilistic or stochastic context-free grammars (PCFGs)

- G = (T, N, S, R, P)
  - T is set of terminals

  - N is set of nonterminals

  - S is the start symbol (one of the nonterminals)

  - R is rules/productions of the form $X \rightarrow \gamma$

  - P(R) gives the probability of each rule.

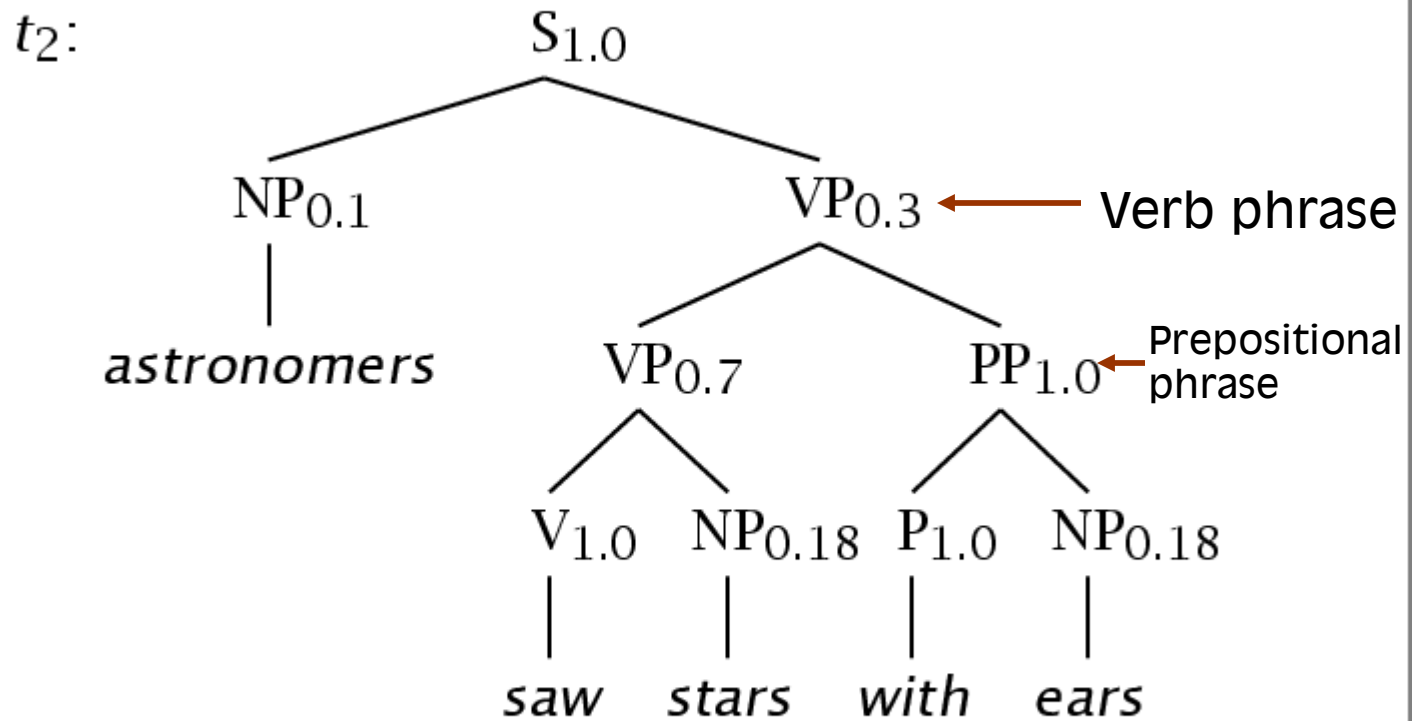$$\forall X \in N, \sum_{X \rightarrow \gamma \in R} P(X \rightarrow \gamma) = 1$$

# A Simple PCFG
# (in Chomsky Normal Form)

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| S | → | NP VP | 1.0 | NP | → | NP PP | 0.4 |
| VP | → | V NP | 0.7 | NP | → | *astronomers* | 0.1 |
| VP | → | VP PP | 0.3 | NP | → | *ears* | 0.18 |
| PP | → | P NP | 1.0 | NP | → | *saw* | 0.04 |
| P | → | *with* | 1.0 | NP | → | *stars* | 0.18 |
| V | → | *saw* | 1.0 | NP | → | *telescope* | 0.1 |

# Parse tree 1:

$t_1$:

$S_{1.0}$

Noun phrase → $NP_{0.1}$

$NP_{0.1}$ — astronomers

$VP_{0.7}$ ← Verb phrase

$V_{1.0}$ — saw

$NP_{0.4}$

$NP_{0.18}$ — stars

$PP_{1.0}$

$P_{1.0}$ — with

$NP_{0.18}$ — ears

# Parse tree 2:

$t_2$: $S_{1.0}$

$NP_{0.1}$ $VP_{0.3}$ ← Verb phrase

astronomers

$VP_{0.7}$ $PP_{1.0}$ ← Prepositional phrase

$V_{1.0}$ $NP_{0.18}$ $P_{1.0}$ $NP_{0.18}$

saw   stars   with   ears

# The 3 questions of PCFG's

- What is the probability of a sentence $w_{1m}$ to a grammar $G$: $P(w_{1m} | G)$ ?

- What is the most likely parse for a sentence
  $\arg\max_t P(t | w_{1m}, G)$ ?

- How can we choose rule probabilities for the grammar $G$ that maximize the probability of a sentence $\arg\max_G P(w_{1m} | G)$ ?
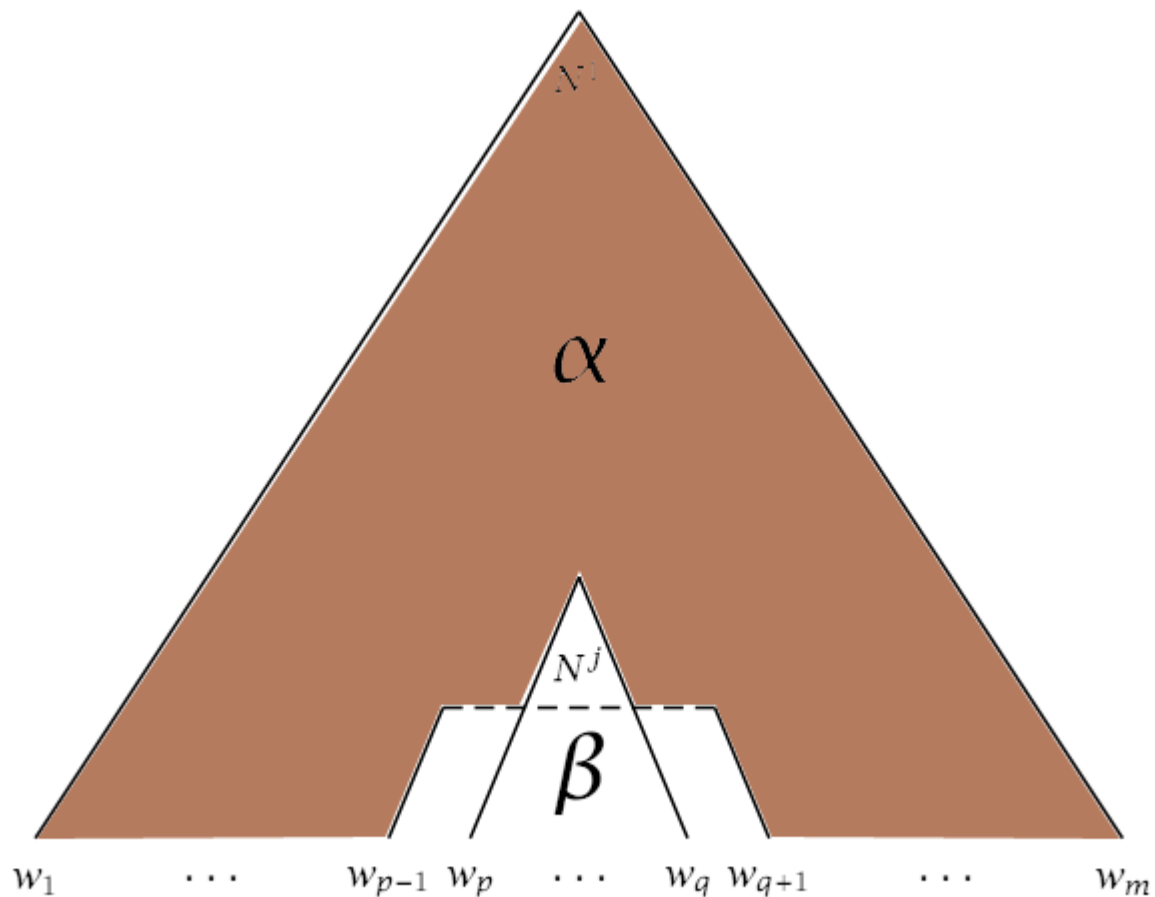
Figure 11.3    Inside and outside probabilities in PCFGs.

Outside probability    $\alpha_j(p,q) = P(w_{1(p-1)}, N^j_{pq} w_{(q+1)} \mid G)$

Inside probability    $\beta_j(p,q) = P(w_{pq} \mid N^j_{pq}, G)$

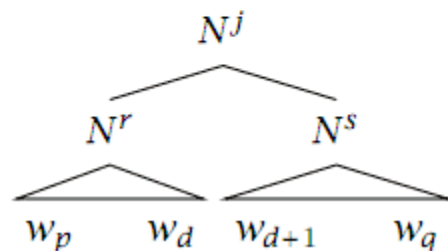*Probabilistic Context Free Grammars*    10

Figure 11.3 Inside and outside probabilities in PCFGs.

Outside probability $\quad \alpha_j(p,q) = P(w_{1(p-1)}, N^j_{pq} w_{(q+1)} \mid G)$

Inside probability $\quad \beta_j(p,q) = P(w_{pq} \mid N^j_{pq}, G)$

*Probabilistic Context Free Grammars*

# The probability of a String

□ Inside algorithm:

$$P(w_{1m}|G) = P(N^1 \overset{*}{\Rightarrow} w_{1m}|G)$$
$$= P(w_{1m}|N^1_{1m}, G) = \beta_1(1, m)$$

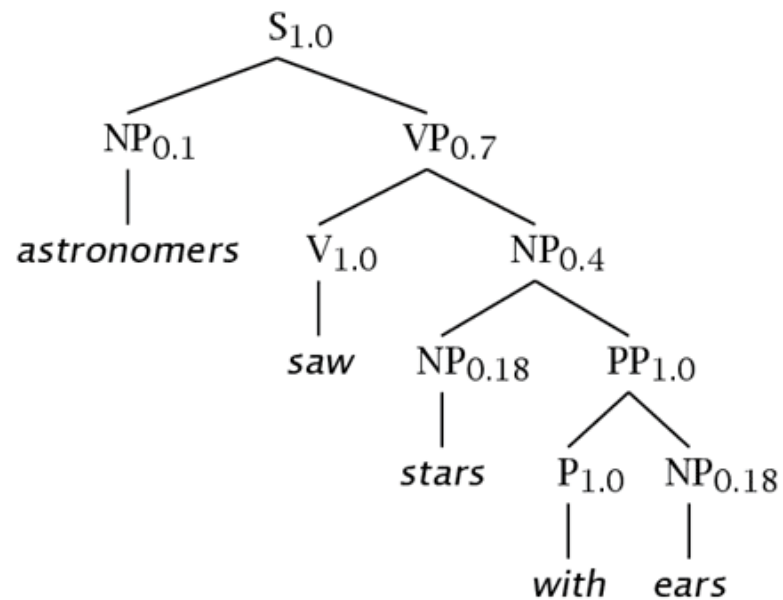**Base case:** We want to find $\beta_j(k, k)$ (the probability of a rule $N^j \to w_k$):

$$\beta_j(k, k) = P(w_k|N^j_{kk}, G)$$
$$= P(N^j \to w_k|G)$$

**Induction:**

# The probability of a String

| | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| 1 | $\beta_{NP}$ = 0.1 | | $\beta_S$ = 0.0126 | | $\beta_S$ = 0.0015876 |
| 2 | | $\beta_{NP}$ = 0.04 $\beta_V$ = 1.0 | $\beta_{VP}$ = 0.126 | | $\beta_{VP}$ = 0.015876 |
| 3 | | | $\beta_{NP}$ = 0.18 | | $\beta_{NP}$ = 0.01296 |
| 4 | | | | $\beta_P$ = 1.0 | $\beta_{PP}$ = 0.18 |
| 5 | | | | | $\beta_{NP}$ = 0.18 |
| | astronomers | saw | stars | with | ears |

# Using outside probabilities

For any $k$, $1 \le k \le m$,

$$P(w_{1m}|G) = \sum_j P(w_{1(k-1)}, w_k, w_{(k+1)m}, N^j_{kk}|G)$$

$$= \sum_j P(w_{1(k-1)}, N^j_{kk}, w_{(k+1)m}|G)$$

$$\times P(w_k|w_{1(k-1)}, N^j_{kk}, w_{(k+1)n}, G)$$

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

**Base Case:** The base case is the probability of the root of the tree being nonterminal $N^i$ with nothing outside it:
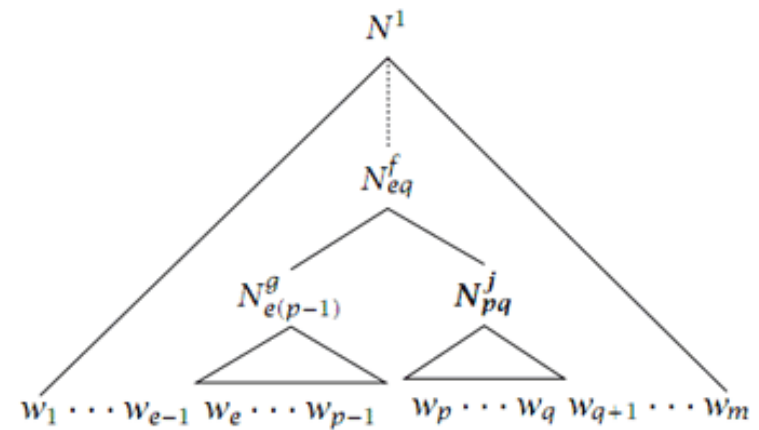
$$\alpha_1(1,m) = 1$$
$$\alpha_j(1,m) = 0 \quad \text{for } j \ne 1$$

# Using outside probabilities



Inductive case: a node

$N_{pq}^j$ might be on the left:

or right branch of the parent node:

# Using outside probabilities

$$\alpha_j(p,q) = \Big[ \sum_{f,g \neq j} \sum_{e=q+1}^{m} P(w_{1(p-1)}, w_{(q+1)m}, N_{pe}^f, N_{pq}^j, N_{(q+1)e}^g) \Big]$$

$$+ \Big[ \sum_{f,g} \sum_{e=1}^{p-1} P(w_{1(p-1)}, w_{(q+1)m}, N_{eq}^f, N_{e(p-1)}^g, N_{pq}^j) \Big]$$

$$= \Big[ \sum_{f,g \neq j} \sum_{e=q+1}^{m} P(w_{1(p-1)}, w_{(e+1)m}, N_{pe}^f) P(N_{pq}^j, N_{(q+1)e}^g | N_{pe}^f)$$

$$\times P(w_{(q+1)e} | N_{(q+1)e}^g) \Big] + \Big[ \sum_{f,g} \sum_{e=1}^{p-1} P(w_{1(e-1)}, w_{(q+1)m}, N_{eq}^f)$$

$$\times P(N_{e(p-1)}^g, N_{pq}^j | N_{eq}^f) P(w_{e(p-1)} | N_{e(p-1)}^g) \Big]$$

$$= \Big[ \sum_{f,g \neq j} \sum_{e=q+1}^{m} \alpha_f(p,e) P(N^f \to N^j N^g) \beta_g(q+1,e) \Big]$$

$$+ \Big[ \sum_{f,g} \sum_{e=1}^{p-1} \alpha_f(e,q) P(N^f \to N^g N^j) \beta_g(e,p-1) \Big]$$

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

$$= P(w_{1m}, N_{pq}^j | G)$$

$$P(w_{1m}, N_{pq} | G) = \sum_{j} \alpha_j(p,q) \beta_j(p,q)$$

# Finding the most likely parse for a sentence

$\delta_i(p,q)$ = the highest inside probability parse of a subtree $N^i_{pq}$

1. Initialization

$$\delta_i(p,p) = P(N^i \rightarrow w_p)$$

2. Induction

$$\delta_i(p,q) = \max_{\substack{1 \le j,k \le n \\ p \le r < q}} P(N^i \rightarrow N^j \; N^k)\delta_j(p,r)\delta_k(r+1,q)$$

Store backtrace

$$\psi_i(p,q) = \underset{(j,k,r)}{\arg\max} \, P(N^i \rightarrow N^j \; N^k)\delta_j(p,r)\delta_k(r+1,q)$$

3. Termination and path readout (by backtracking).

$$P(\hat{t}) = \delta_1(1,m)$$

*Probabilistic Context Free Grammars*                                   17

# Training a PCFG with the Inside-Outside algorithm

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

$$\alpha_j(p,q)\beta_j(p,q) = P(N^1 \overset{*}{\Rightarrow} w_{1m}, N^j \overset{*}{\Rightarrow} w_{pq}|G)$$

$$= P(N^1 \overset{*}{\Rightarrow} w_{1m}|G)P(N^j \overset{*}{\Rightarrow} w_{pq}|N^1 \overset{*}{\Rightarrow} w_{1m},G)$$

$$\pi = P(N^1 \overset{*}{\Rightarrow} w_{1m})$$

$$P(N^j \overset{*}{\Rightarrow} w_{pq}|N^1 \overset{*}{\Rightarrow} w_{1m},G) = \frac{\alpha_j(p,q)\beta_j(p,q)}{\pi}$$

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

$\forall r,s,p < q:$

$$P(N^j \rightarrow N^r N^s \overset{*}{\Rightarrow} w_{pq}|N^1 \overset{*}{\Rightarrow} w_{1m},G)$$

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

*Probabilistic Context Free Grammars*

18

# Training a PCFG with the Inside-Outside algorithm

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

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

Now for the maximization step, we want:

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

So, the reestimation formula is:

$\hat{P}(N^j \rightarrow N^r N^s) =$

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

Similarly for preterminals,

$$P(N^j \rightarrow w^k | N^1 \overset{*}{\Rightarrow} w_{1m}, G) = \frac{\sum_{h=1}^{m} \alpha_j(h,h) P(N^j \rightarrow w_h, w_h = w^k)}{\pi}$$

$$= \frac{\sum_{h=1}^{m} \alpha_j(h,h) P(w_h = w^k) \beta_j(h,h)}{\pi}$$

$$\hat{P}(N^j \rightarrow w^k) = \frac{\sum_{h=1}^{m} \alpha_j(h,h) P(w_h = w^k) \beta_j(h,h)}{\sum_{p=1}^{m} \sum_{q=p}^{m} \alpha_j(p,q) \beta_j(p,q)}$$

*Probabilistic Context Free Grammars*

# Training a PCFG with the Inside-Outside algorithm

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

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

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

$$\hat{P}(N^j \rightarrow N^r N^s) = \frac{\sum_{i=1}^{\omega} \sum_{p=1}^{m_i-1} \sum_{q=p+1}^{m_i} f_i(p,q,j,r,s)}{\sum_{i=1}^{\omega} \sum_{p=1}^{m_i} \sum_{q=p}^{m_i} h_i(p,q,j)} \quad \leftarrow \text{Estimation}$$

and

$$\hat{P}(N^j \rightarrow w^k) = \frac{\sum_{i=1}^{\omega} \sum_{h=1}^{m_i} g_i(h,j,k)}{\sum_{i=1}^{\omega} \sum_{p=1}^{m_i} \sum_{q=p}^{m_i} h_i(p,q,j)} \quad \leftarrow \text{Reestimation}$$

$$P(W|G_{i+1}) \geq P(W|G_i).$$

# Problems with
# the Inside-Outside algorithm

- Training PCFGs is way slow: For each sentence, each iteration of training is $O(m^3n^3)$, being $m$ the length of the sentence and $n$ the number of non-terminals in the grammar.

- Algorithm is very sensitive to initialization parameters.

# Some features of PCFGs

- Better than grammars based on HMMs, as the diversity and size of a corpus of texts expands.

- PCFG's don't take lexical context into account.

- For that reason, it does not give a plausibility of different parses.

- Robustness: Real text has grammatical errors. Just give implausible sentences a low probability

- PCFG's have certain biases: a smaller tree has more probability than a larger tree.

# Acknowledgments

- Some material of this presentation was taken from Lecture 2 on Statistical Parsing by Christopher Manning, as well some graphics and examples from Chapter 11 of the book *Foundations of Statistical Natural Language Processing*  by Christopher D. Manning and Hinrich Schütze