

context-free grammar $G = (N, \Sigma, S, P)$

- N finite set (nonterminals)
- Σ finite set (terminals) $N \cap \Sigma = \emptyset$
- $S \in N$ (initial nonterminal)
- P finite set (rules)

$$A \rightarrow \alpha \quad A \in N, \alpha \in (N \cup \Sigma)^*$$

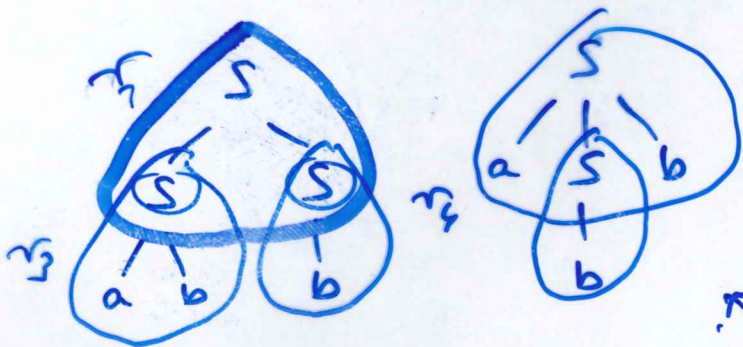
Example:

- $r_1: S \rightarrow SS$
- $r_2: S \rightarrow aSb$
- $r_3: S \rightarrow ab$
- $r_4: S \rightarrow b$

$$S \rightarrow A \underbrace{S B}_{S'}$$

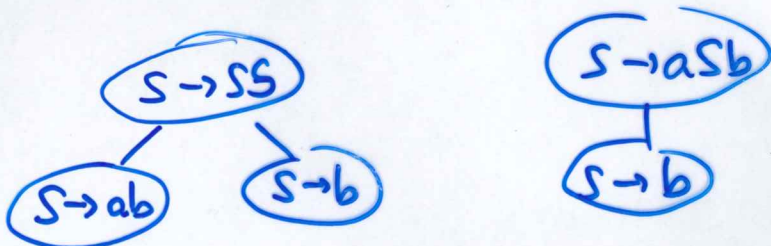
$$S \Rightarrow AS' \Rightarrow ASB$$

parse trees for abb:



$r_1 r_3 r_4$
 bzw. $r_2 r_4$

abstract syntax trees:



probabilistic cfg (G, p)

- $G = (N, \Sigma, S, P)$ cfg

- $p: P \rightarrow [0, 1]$

p is proper if $\forall A \in N: \sum_{\alpha = (A \rightarrow \alpha) \text{ in } P} p(\alpha) = 1$

probability of parse tree $d: d = r_1 \dots r_n$

if $d = r_1 \dots r_n$, then $P(r_1 \dots r_n) = p(r_1) \dots p(r_n)$

probability of sentence $w \in \Sigma^*$: $P(r_1 \dots r_n | (G, p))$

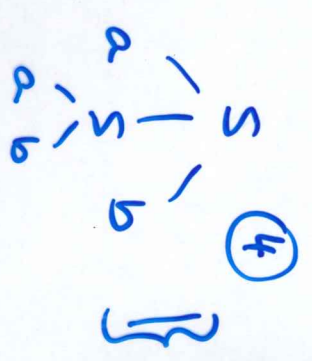
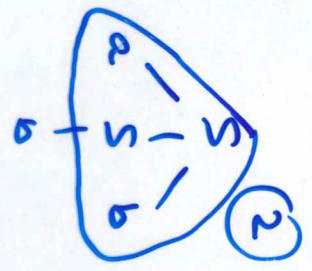
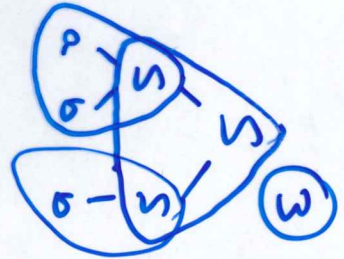
$$P(w | (G, p)) = \sum_{\substack{d \text{ parse tree} \\ \text{of } w}} P(d | (G, p))$$

\approx max
 \uparrow

grammar induction

given: Corpus

$C = \{$



multiset

readoff cfg:

$S \rightarrow SS, S \rightarrow ab, S \rightarrow b, S \rightarrow aSb$

frequency analysis:

r	$S \rightarrow SS$	$S \rightarrow ab$	$S \rightarrow b$	$S \rightarrow aSb$
number of occ. of r	3	7	5	6

sum of number of occ. of all rules with S

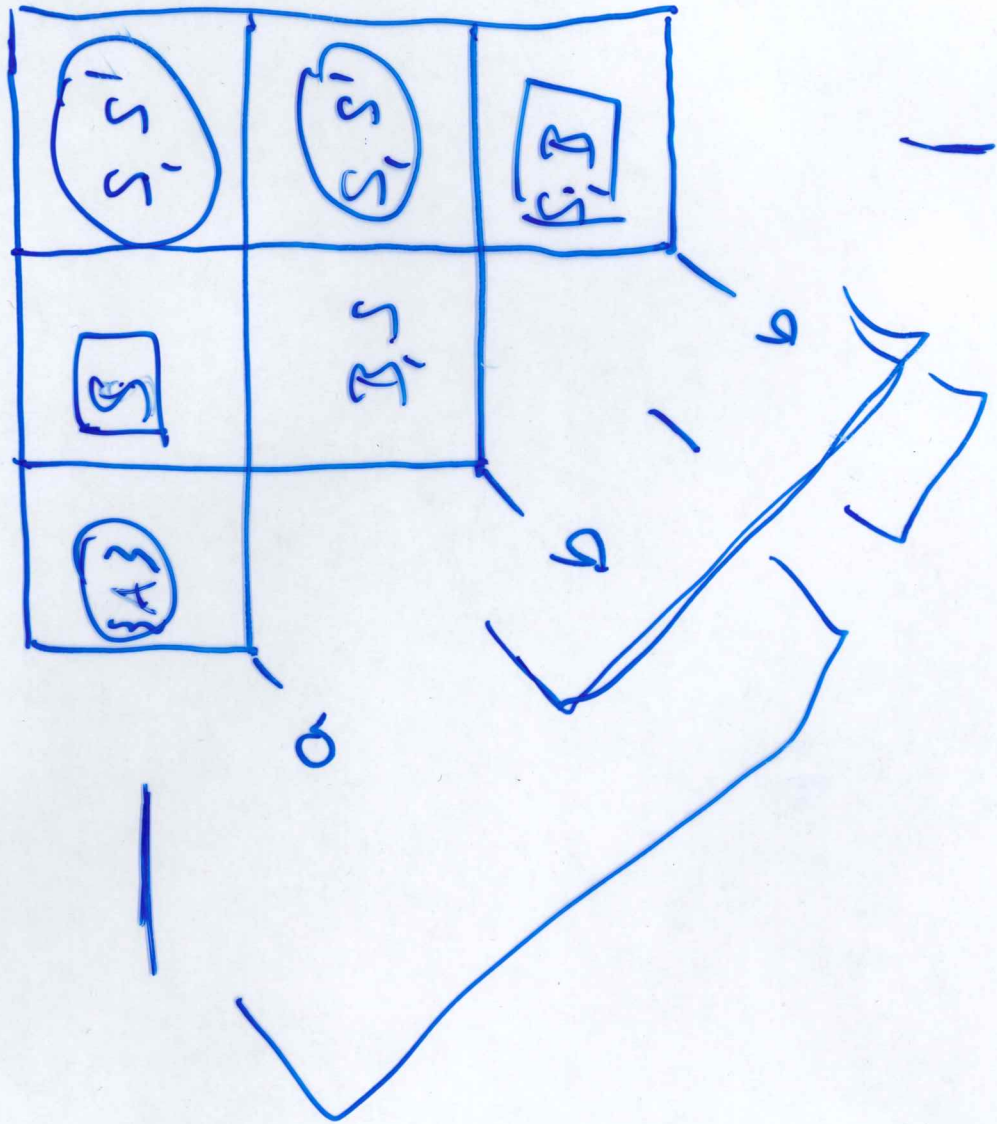
21

$P(r)$	$\frac{3}{21} \approx 0.143$	$\frac{7}{21} \approx 0.33$	$\frac{5}{21} \approx 0.238$	$\frac{6}{21} \approx 0.286$
--------	------------------------------	-----------------------------	------------------------------	------------------------------

CYK - algorithm (Cocke - Younger - Karami)

Chomsky NF

- $S \rightarrow SS$
- $S \rightarrow AS'$
- $S' \rightarrow SB$
- $S \rightarrow AB$
- $S \rightarrow b$
- $A \rightarrow a$
- $B \rightarrow b$



An adaptation of Knuth's algorithm [Knu77]

Require: PCFG (\mathcal{G}, p) with $\mathcal{G} = (N, \Sigma, S, R)$ and $p(r) > 0$ for every r

Ensure: $\hat{d} = \arg \max_{d \in D_{\mathcal{G}}} P(d)$

▷ family $\delta = (\delta_A \mid A \in N)$ with $\delta_A \in (D_{\mathcal{G}}(A, \Sigma^*) \times [0, 1]) \cup \{(\perp, 0)\}$, $U, Q \subseteq N$

- 1: **for all** $A \in N$ **do**
- 2: $\delta_A \leftarrow (\perp, 0)$;
- 3: **for all** $r = (A \rightarrow u) \in R$ with $u \in \Sigma^*$ **do**
- 4: **if** $p(r) > (\delta_A)_2$ **then**
- 5: $\delta_A \leftarrow (r, p(r))$
- 6: $Q \leftarrow \emptyset$;
- 7: $U \leftarrow N$;
- 8: **while** $U \neq \emptyset$ **do**
- 9: $A \leftarrow \arg \max_{B \in U} (\delta_B)_2$;
- 10: $U \leftarrow U \setminus \{A\}$;
- 11: $Q \leftarrow Q \cup \{A\}$;
- 12: **for all** $r = (B \rightarrow u_0 B_1 u_1 \dots B_k u_k) \in R$ **do**
- 13: **if** $B \in U$ and $A \in \{B_1, \dots, B_k\} \subseteq Q$ **then**
- 14: $d := r((\delta_{B_1})_1, \dots, (\delta_{B_k})_1)$;
- 15: $s := p(r) \cdot \prod_{i=1}^k (\delta_{B_i})_2$;
- 16: **if** $s \geq (\delta_B)_2$ **then**
- 17: $\delta_B \leftarrow (d, s)$
- 18: $\hat{d} \leftarrow (\delta_S)_1$

▷ derivation viewed as abstract syntax tree