

CS 378 Lecture 12

Constituency Parsing

Announcements

- Midterm: Weds - Sat of next week
- A3 due Tues
- A2 back soon

Recap

Viterbi: dynamic programming alg to
compute $\max_{\bar{y}} P(\bar{y}, \bar{x})$ for an HMM

(similar alg today for parsing: CKY)

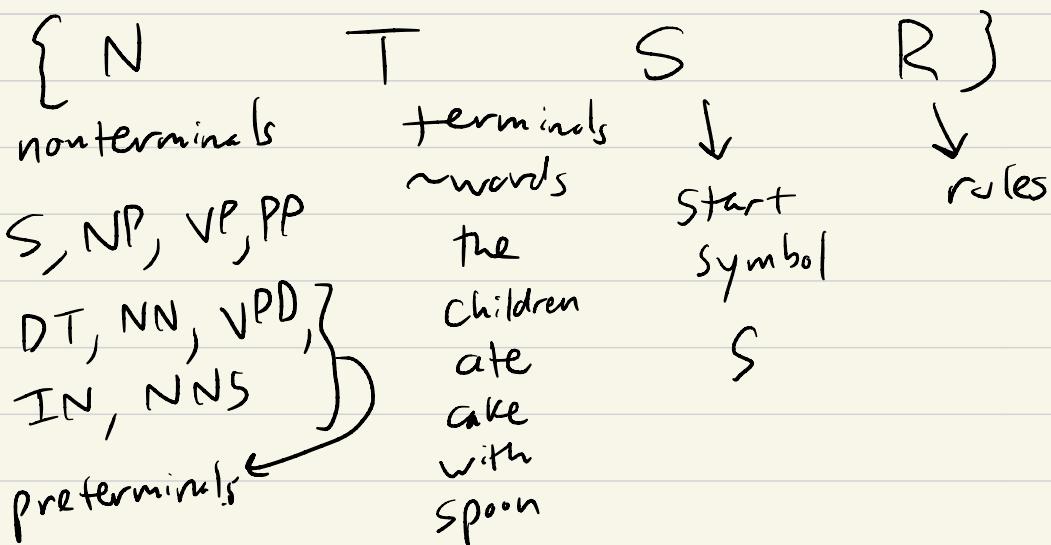
Beam search: approximates Viterbi by
only preserving K states at a timestep

Today Constituency intro

Probabilistic context-free grammars
CKY algorithm

Constituency see slides

Context-free grammars



binary

$S \rightarrow NP \ VP \ 1.0$

$VP \rightarrow VBD \ NP \ 1.0$

$NP \rightarrow DT \ NN \ 0.5$

$NP \rightarrow DT \ NNS \ 0.5$

unary

$DT \rightarrow \text{the} \ 1.0$

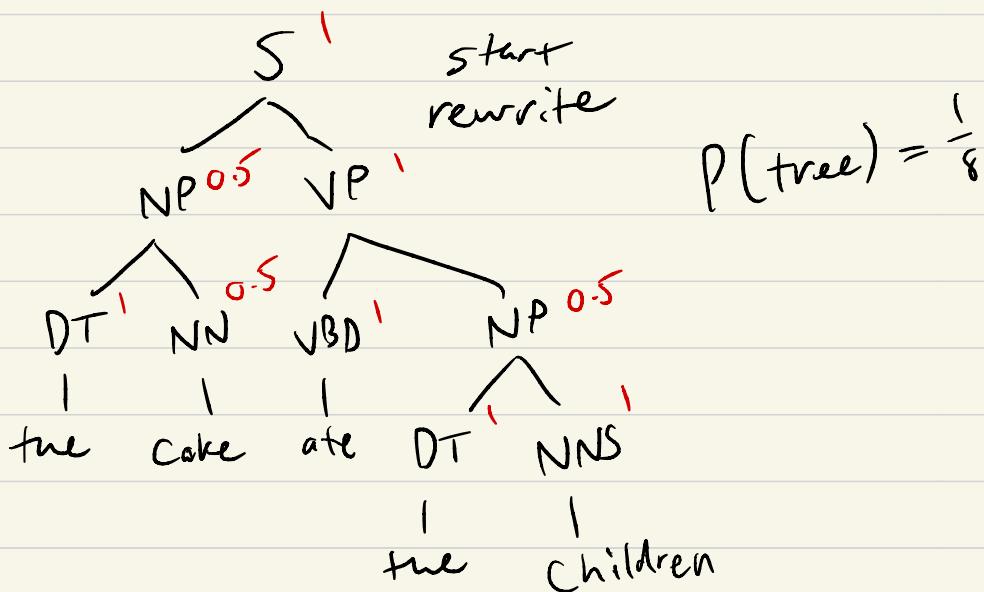
$NN \rightarrow \text{children} \ 1.0$

$NN \rightarrow \text{cake} \ 0.5$

$NN \rightarrow \text{spoon} \ 0.5$

$VBD \rightarrow \text{ate} \ 1.0$

CFG defines a set of trees



Probabilistic CFGs

Each rule has a prob.
Probs normalize per parent

$$P(\text{rule} \mid \text{NP}) = \begin{cases} \text{NP} \rightarrow \text{DT NN} & 0.5 \\ \text{NP} \rightarrow \text{DT NNS} & 0.5 \end{cases}$$

(~ HMM transitions)

$$P(T) = \prod_{\text{rules}} P(\text{rule} \mid \text{parent(rule)})$$

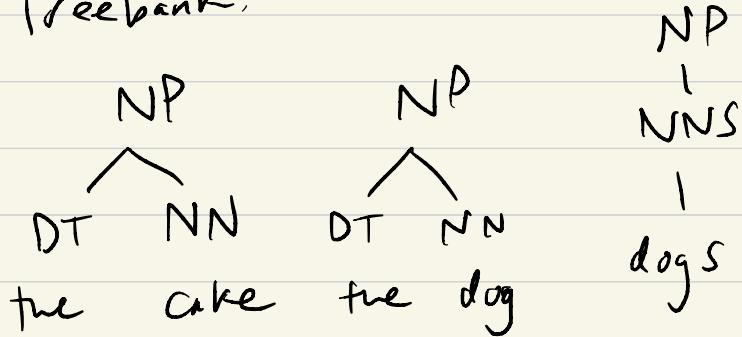
Building a parser

- Input: tree bank (sents w/ labeled trees)
Output: model + way to compute $\arg\max_T P(T \mid \bar{x})$
- ① Grammar preprocessing (binarization) next time
 - ② Computing grammar + getting probs
 - ③ Parsing

① Make every tree binary

② Read off grammar + count rules

Treebank:



$$P(w|NN) = \begin{cases} 1/2 & \text{cake} \\ 1/2 & \text{dog} \end{cases}$$

$$P(w|NNS) = \{1.0 \text{ dogs}\}$$

$$P(r|NP) = \begin{cases} 2/3 & \text{DT } NN \\ 1/3 & \text{NNS} \end{cases}$$

This is the maximum likelihood estimate
of params for this data

③ CKY algorithm

Input: PCFG

sentence \bar{x}

Output: $\underset{T}{\operatorname{argmax}} P(T | \bar{x})$

most likely tree for
 \bar{x}

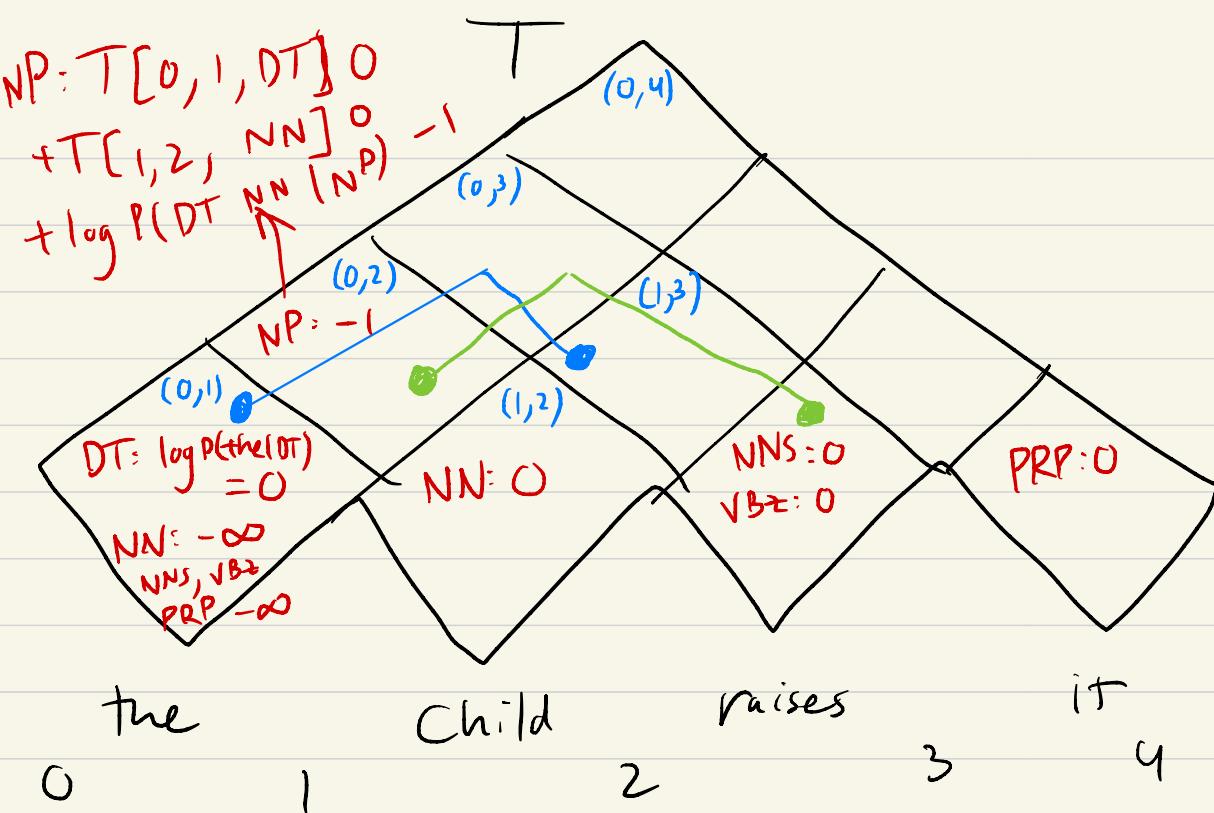
$$= \underset{T}{\operatorname{argmax}} P(T, \bar{x})$$

dynamic program: track the best score
for each nonterminal over each
span of the sentence

$T(i, j, X) = \text{score} (\sim \log \text{prob}) \text{ of the}$
 $\text{best way to build symbol}$
 $X \text{ over the span } (i, j)$

sent len
+ sent len
+ number of symbols

\uparrow
python list indices



Grammar:

$$\begin{array}{ll}
 DT \rightarrow \text{the} & | \\
 NN \rightarrow \text{child} & | \\
 NNS \rightarrow \text{raises} & | \\
 VBZ \rightarrow \text{raises} & | \\
 PRP \rightarrow \text{it} & |
 \end{array}
 \qquad
 \begin{array}{ll}
 S \rightarrow NP \quad VP & | \\
 NP \rightarrow DT \quad NN & |_{1/2} \\
 NP \rightarrow NN \quad NNS & |_{1/2} \\
 VP \rightarrow VBZ \quad PRP & |
 \end{array}$$

$$\log(1/2) = -1$$

CKY

base: $T(i, i+1, X) = \log P(w_i | X)$

recursive:

$$T(i, j, X) = \max_{\substack{k: i < k < j \\ \text{split point}}} \max_{\substack{X \rightarrow X_1 X_2 \\ \text{rule}}} \left[\log P(X \rightarrow X_1 X_2) + T[i, k, X_1] + T[k, j, X_2] \right]$$



cell (0,3)

(0,1) + (1,3)

(0,2) + (2,3)