

Probabilistic Context-Free Grammars

CFGs $\{N, T, S, R\}$

nonterminals, terminals, start, rules

S, NP, VP, PP

the, children,
ate, cake,
with,
spoon

S binary

Unary
 $VP \rightarrow VBD$ $3/4$

DT, NN, VBD, IN, NNS

POS: preterminals

$1 \quad S \rightarrow NP \quad VP$

$1/4 \quad VP \rightarrow VBD \quad NP$

$1/2 \quad NP \rightarrow DT \quad NN$

$1/2 \quad NP \rightarrow DT \quad NNS$

$DT \rightarrow the$
 $NNS \rightarrow children$

$NN \rightarrow cake$

$NN \rightarrow spoon$

$VBD \rightarrow ate$

PCFG: rules have probs, probs sum to one

$P(\text{rule} \mid \text{parent}(\text{rule}))$ per parent

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

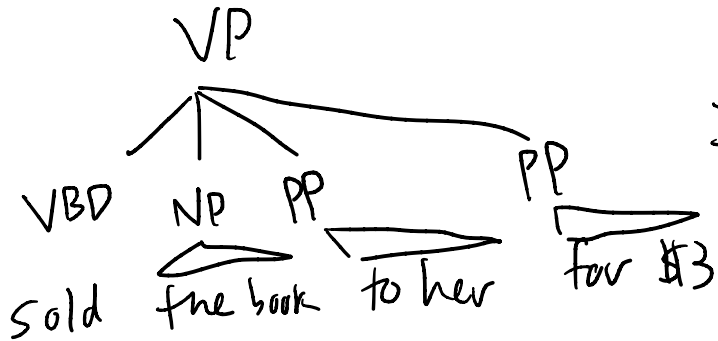
Treebank of sents labeled with trees

- ## ① Binarization

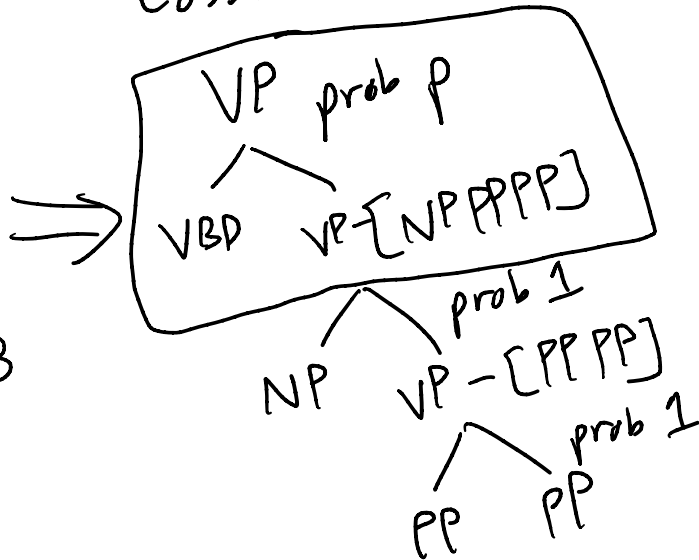
- ② Estimation of rule probs
(count + normalizing)
(maximum likelihood est.)

- ③ Inference
 $\operatorname{argmax}_T P(T|x)$

Binarization



Lossless



Lossy

