

# Probabilistic Context-Free Grammars

Michael Collins, Columbia University

# Overview

- ▶ Probabilistic Context-Free Grammars (PCFGs)
- ▶ The CKY Algorithm for parsing with PCFGs

# A Probabilistic Context-Free Grammar (PCFG)

S	⇒	NP	VP	1.0
VP	⇒	Vi		0.4
VP	⇒	Vt	NP	0.4
VP	⇒	VP	PP	0.2
NP	⇒	DT	NN	0.3
NP	⇒	NP	PP	0.7
PP	⇒	P	NP	1.0

Vi	⇒	sleeps	1.0
Vt	⇒	saw	1.0
NN	⇒	man	0.7
NN	⇒	woman	0.2
NN	⇒	telescope	0.1
DT	⇒	the	1.0
IN	⇒	with	0.5
IN	⇒	in	0.5

- Probability of a tree  $t$  with rules

$$\alpha_1 \rightarrow \beta_1, \alpha_2 \rightarrow \beta_2, \dots, \alpha_n \rightarrow \beta_n$$

is  $p(t) = \prod_{i=1}^n q(\alpha_i \rightarrow \beta_i)$  where  $q(\alpha \rightarrow \beta)$  is the probability for rule  $\alpha \rightarrow \beta$ .

DERIVATION

S

RULES USED

PROBABILITY

DERIVATION

S

NP VP

RULES USED

$S \rightarrow NP VP$

PROBABILITY

1.0

## DERIVATION

S

NP VP

DT NN VP

## RULES USED

$S \rightarrow NP VP$

$NP \rightarrow DT NN$

## PROBABILITY

1.0

0.3

## DERIVATION

S

NP VP

DT NN VP

the NN VP

## RULES USED

$S \rightarrow NP VP$

$NP \rightarrow DT NN$

$DT \rightarrow \text{the}$

## PROBABILITY

1.0

0.3

1.0

## DERIVATION

S

NP VP

DT NN VP

the NN VP

the dog VP

## RULES USED

$S \rightarrow NP VP$

$NP \rightarrow DT NN$

$DT \rightarrow \text{the}$

$NN \rightarrow \text{dog}$

## PROBABILITY

1.0

0.3

1.0

0.1



**DERIVATION**

S

NP VP

DT NN VP

the NN VP

the dog VP

the dog Vi

**RULES USED** $S \rightarrow NP VP$  $NP \rightarrow DT NN$  $DT \rightarrow \text{the}$  $NN \rightarrow \text{dog}$  $VP \rightarrow V_i$ **PROBABILITY**

1.0

0.3

1.0

0.1

0.4

**DERIVATION**

S

NP VP

DT NN VP

the NN VP

the dog VP

the dog Vi

the dog laughs

**RULES USED** $S \rightarrow NP VP$  $NP \rightarrow DT NN$  $DT \rightarrow \text{the}$  $NN \rightarrow \text{dog}$  $VP \rightarrow V_i$  $V_i \rightarrow \text{laughs}$ **PROBABILITY**

1.0

0.3

1.0

0.1

0.4

0.5

# Properties of PCFGs

- ▶ Assigns a probability to each *left-most derivation*, or parse-tree, allowed by the underlying CFG

# Properties of PCFGs

- ▶ Assigns a probability to each *left-most derivation*, or parse-tree, allowed by the underlying CFG
- ▶ Say we have a sentence  $s$ , set of derivations for that sentence is  $\mathcal{T}(s)$ . Then a PCFG assigns a probability  $p(t)$  to each member of  $\mathcal{T}(s)$ . i.e., *we now have a ranking in order of probability.*

# Properties of PCFGs

- ▶ Assigns a probability to each *left-most derivation*, or parse-tree, allowed by the underlying CFG
- ▶ Say we have a sentence  $s$ , set of derivations for that sentence is  $\mathcal{T}(s)$ . Then a PCFG assigns a probability  $p(t)$  to each member of  $\mathcal{T}(s)$ . i.e., *we now have a ranking in order of probability*.
- ▶ The most likely parse tree for a sentence  $s$  is

$$\arg \max_{t \in \mathcal{T}(s)} p(t)$$



## Deriving a PCFG from a Treebank

- ▶ Given a set of example trees (a treebank), the underlying CFG can simply be **all rules seen in the corpus**
- ▶ Maximum Likelihood estimates:

$$q_{ML}(\alpha \rightarrow \beta) = \frac{\text{Count}(\alpha \rightarrow \beta)}{\text{Count}(\alpha)}$$

where the counts are taken from a training set of example trees.

- ▶ **If the training data is generated by a PCFG**, then as the training data size goes to infinity, the maximum-likelihood PCFG will converge to the same distribution as the “true” PCFG.

# PCFGs

Booth and Thompson (1973) showed that a CFG with rule probabilities correctly defines a distribution over the set of derivations provided that:

1. The rule probabilities define conditional distributions over the different ways of rewriting each non-terminal.
2. A technical condition on the rule probabilities ensuring that the probability of the derivation terminating in a finite number of steps is 1. (This condition is not really a practical concern.)



# Parsing with a PCFG

- ▶ Given a PCFG and a sentence  $s$ , define  $\mathcal{T}(s)$  to be the set of trees with  $s$  as the yield.
- ▶ Given a PCFG and a sentence  $s$ , how do we find

$$\arg \max_{t \in \mathcal{T}(s)} p(t)$$

# Chomsky Normal Form

A context free grammar  $G = (N, \Sigma, R, S)$  in Chomsky Normal Form is as follows

- ▶  $N$  is a set of non-terminal symbols
- ▶  $\Sigma$  is a set of terminal symbols
- ▶  $R$  is a set of rules which take one of two forms:
  - ▶  $X \rightarrow Y_1Y_2$  for  $X \in N$ , and  $Y_1, Y_2 \in N$
  - ▶  $X \rightarrow Y$  for  $X \in N$ , and  $Y \in \Sigma$
- ▶  $S \in N$  is a distinguished start symbol

# A Dynamic Programming Algorithm

- ▶ Given a PCFG and a sentence  $s$ , how do we find

$$\max_{t \in \mathcal{T}(s)} p(t)$$

- ▶ Notation:

$n$  = number of words in the sentence

$w_i$  =  $i$ 'th word in the sentence

$N$  = the set of non-terminals in the grammar

$S$  = the start symbol in the grammar

- ▶ Define a dynamic programming table

$\pi[i, j, X]$  = maximum probability of a constituent with non-terminal  $X$   
spanning words  $i \dots j$  inclusive

- ▶ Our goal is to calculate  $\max_{t \in \mathcal{T}(s)} p(t) = \pi[1, n, S]$

## An Example

the dog saw the man with the telescope

# A Dynamic Programming Algorithm

- ▶ Base case definition: for all  $i = 1 \dots n$ , for  $X \in N$

$$\pi[i, i, X] = q(X \rightarrow w_i)$$

(note: define  $q(X \rightarrow w_i) = 0$  if  $X \rightarrow w_i$  is not in the grammar)

- ▶ Recursive definition: for all  $i = 1 \dots n$ ,  $j = (i + 1) \dots n$ ,  $X \in N$ ,

$$\pi(i, j, X) = \max_{\substack{X \rightarrow YZ \in R, \\ s \in \{i \dots (j-1)\}}} (q(X \rightarrow YZ) \times \pi(i, s, Y) \times \pi(s + 1, j, Z))$$

## An Example

$$\pi(i, j, X) = \max_{\substack{X \rightarrow YZ \in R, \\ s \in \{i \dots (j-1)\}}} (q(X \rightarrow YZ) \times \pi(i, s, Y) \times \pi(s + 1, j, Z))$$

the dog saw the man with the telescope

# The Full Dynamic Programming Algorithm

**Input:** a sentence  $s = x_1 \dots x_n$ , a PCFG  $G = (N, \Sigma, S, R, q)$ .

**Initialization:**

For all  $i \in \{1 \dots n\}$ , for all  $X \in N$ ,

$$\pi(i, i, X) = \begin{cases} q(X \rightarrow x_i) & \text{if } X \rightarrow x_i \in R \\ 0 & \text{otherwise} \end{cases}$$

**Algorithm:**

- ▶ For  $l = 1 \dots (n - 1)$ 
  - ▶ For  $i = 1 \dots (n - l)$ 
    - ▶ Set  $j = i + l$
    - ▶ For all  $X \in N$ , calculate

$$\pi(i, j, X) = \max_{\substack{X \rightarrow YZ \in R, \\ s \in \{i \dots (j-1)\}}} (q(X \rightarrow YZ) \times \pi(i, s, Y) \times \pi(s + 1, j, Z))$$

and

$$bp(i, j, X) = \arg \max_{\substack{X \rightarrow YZ \in R, \\ s \in \{i \dots (j-1)\}}} (q(X \rightarrow YZ) \times \pi(i, s, Y) \times \pi(s + 1, j, Z))$$

# A Dynamic Programming Algorithm for the Sum

- ▶ Given a PCFG and a sentence  $s$ , how do we find

$$\sum_{t \in \mathcal{T}(s)} p(t)$$

- ▶ Notation:

$n$  = number of words in the sentence

$w_i$  =  $i$ 'th word in the sentence

$N$  = the set of non-terminals in the grammar

$S$  = the start symbol in the grammar

- ▶ Define a dynamic programming table

$\pi[i, j, X]$  = sum of probabilities for constituent with non-terminal  $X$   
spanning words  $i \dots j$  inclusive

- ▶ Our goal is to calculate  $\sum_{t \in \mathcal{T}(s)} p(t) = \pi[1, n, S]$



# Summary

- ▶ PCFGs augments CFGs by including a probability for each rule in the grammar.
- ▶ The probability for a parse tree is the product of probabilities for the rules in the tree
- ▶ To build a PCFG-parsed parser:
  1. Learn a PCFG from a treebank
  2. Given a test data sentence, use the CKY algorithm to compute the highest probability tree for the sentence under the PCFG