

CMPT-413

Computational Linguistics

Anoop Sarkar
<http://www.cs.sfu.ca/~anoop>

March 17, 2011

Probabilistic CFG (PCFG)

<i>S</i>	\rightarrow	<i>NP VP</i>	1
<i>VP</i>	\rightarrow	<i>V NP</i>	0.9
<i>VP</i>	\rightarrow	<i>VP PP</i>	0.1
<i>PP</i>	\rightarrow	<i>P NP</i>	1
<i>NP</i>	\rightarrow	<i>NP PP</i>	0.25
<i>NP</i>	\rightarrow	<i>Calvin</i>	0.25
<i>NP</i>	\rightarrow	<i>monsters</i>	0.25
<i>NP</i>	\rightarrow	<i>school</i>	0.25
<i>V</i>	\rightarrow	<i>imagined</i>	1
<i>P</i>	\rightarrow	<i>in</i>	1

$$P(\text{input}) = \sum_{\text{tree}} P(\text{tree} \mid \text{input})$$

$$P(\text{Calvin imagined monsters in school}) = ?$$

Notice that $P(VP \rightarrow V NP) + P(VP \rightarrow VP PP) = 1.0$

Probabilistic CFG (PCFG)

$P(\textit{Calvin imagined monsters in school}) = ?$

```
(S (NP Calvin)
  (VP (V imagined)
    (NP (NP monsters)
      (PP (P in)
        (NP school))))))
```

```
(S (NP Calvin)
  (VP (VP (V imagined)
    (NP monsters))
    (PP (P in)
      (NP school))))
```

Probabilistic CFG (PCFG)

(S (NP Calvin)
 (VP (V imagined)
 (NP (NP monsters)
 (PP (P in)
 (NP school))))))

$$\begin{aligned}P(\text{tree}_1) &= P(S \rightarrow NP VP) \times P(NP \rightarrow Calvin) \times P(VP \rightarrow V NP) \times \\&\quad P(V \rightarrow imagined) \times P(NP \rightarrow NP PP) \times P(NP \rightarrow monsters) \times \\&\quad P(PP \rightarrow P NP) \times P(P \rightarrow in) \times P(NP \rightarrow school) \\&= 1 \times 0.25 \times 0.9 \times 1 \times 0.25 \times 0.25 \times 1 \times 1 \times 0.25 = .003515625\end{aligned}$$

Probabilistic CFG (PCFG)

(S (NP Calvin)
 (VP (VP (V imagined)
 (NP monsters))
 (P (P in)
 (NP school)))))

$$\begin{aligned}P(\text{tree}_2) &= P(S \rightarrow NP VP) \times P(NP \rightarrow Calvin) \times P(VP \rightarrow VP PP) \times \\&\quad P(VP \rightarrow V NP) \times P(V \rightarrow imagined) \times P(NP \rightarrow monsters) \times \\&\quad P(PP \rightarrow P NP) \times P(P \rightarrow in) \times P(NP \rightarrow school) \\&= 1 \times 0.25 \times 0.1 \times 0.9 \times 1 \times 0.25 \times 1 \times 1 \times 0.25 = .00140625\end{aligned}$$

Probabilistic CFG (PCFG)

$$\begin{aligned}P(\text{Calvin imagined monsters in school}) &= P(\text{tree}_1) + P(\text{tree}_2) \\&= .003515625 + .00140625 \\&= .004921875\end{aligned}$$

$$\text{Most likely tree is } \text{tree}_1 = \arg \max_{\text{tree}} P(\text{tree} \mid \text{input})$$

```
(S (NP Calvin)
  (VP (V imagined)
    (NP (NP monsters)
      (PP (P in)
        (NP school))))))
```

```
(S (NP Calvin)
  (VP (VP (V imagined)
    (NP monsters))
    (PP (P in)
      (NP school))))
```

PCFG

- ▶ Central condition: $\sum_{\alpha} P(A \rightarrow \alpha) = 1$
- ▶ Called a *proper* PCFG if this condition holds
- ▶ Note that this means $P(A \rightarrow \alpha) = P(\alpha \mid A) = \frac{f(A, \alpha)}{f(A)}$
- ▶ $P(T \mid S) = \frac{P(T, S)}{P(S)} = P(T, S) = \prod_i P(RHS_i \mid LHS_i)$

- ▶ What is the PCFG that can be extracted from this single tree:
(S (NP (Det the) (NP man))
 (VP (VP (V played)
 (NP (Det a) (NP game)))
 (PP (P with)
 (NP (Det the) (NP dog))))))
- ▶ How many different rhs α exist for $A \rightarrow \alpha$ where A can be S , NP , VP , PP , Det , N , V , P

PCFG

<i>S</i>	\rightarrow	<i>NP VP</i>	$c = 1$	$p = 1/1$	$= 1.0$
<i>NP</i>	\rightarrow	<i>Det NP</i>	$c = 3$	$p = 3/6$	$= 0.5$
<i>NP</i>	\rightarrow	<i>man</i>	$c = 1$	$p = 1/6$	$= 0.1667$
<i>NP</i>	\rightarrow	<i>game</i>	$c = 1$	$p = 1/6$	$= 0.1667$
<i>NP</i>	\rightarrow	<i>dog</i>	$c = 1$	$p = 1/6$	$= 0.1667$
<i>VP</i>	\rightarrow	<i>VP PP</i>	$c = 1$	$p = 1/2$	$= 0.5$
<i>VP</i>	\rightarrow	<i>V NP</i>	$c = 1$	$p = 1/2$	$= 0.5$
<i>PP</i>	\rightarrow	<i>P NP</i>	$c = 1$	$p = 1/1$	$= 1.0$
<i>Det</i>	\rightarrow	<i>the</i>	$c = 2$	$p = 2/3$	$= 0.67$
<i>Det</i>	\rightarrow	<i>a</i>	$c = 1$	$p = 1/3$	$= 0.33$
<i>V</i>	\rightarrow	<i>played</i>	$c = 1$	$p = 1/1$	$= 1.0$
<i>P</i>	\rightarrow	<i>with</i>	$c = 1$	$p = 1/1$	$= 1.0$

- ▶ We can do this with multiple trees. Simply count occurrences of CFG rules over all the trees.
- ▶ A repository of such trees labelled by a human is called a TreeBank.

Ambiguity

- ▶ Part of Speech ambiguity

saw → noun

saw → verb

- ▶ Structural ambiguity: Prepositional Phrases

I saw (the man) with the telescope

I saw (the man with the telescope)

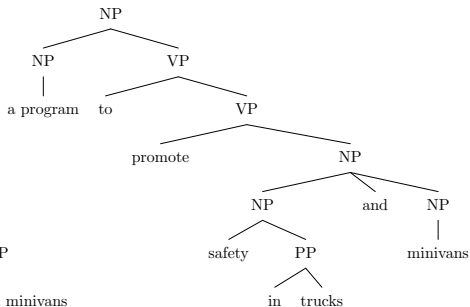
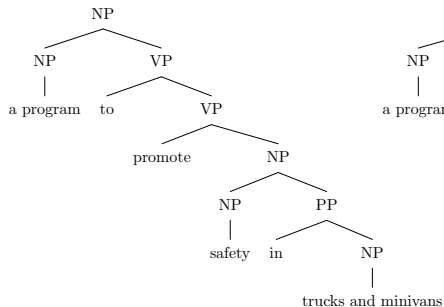
- ▶ Structural ambiguity: Coordination

a program to promote safety in ((trucks) and (minivans))

a program to promote ((safety in trucks) and (minivans))

((a program to promote safety in trucks) and (minivans))

Ambiguity ← attachment choice in alternative parses



Parsing as a machine learning problem

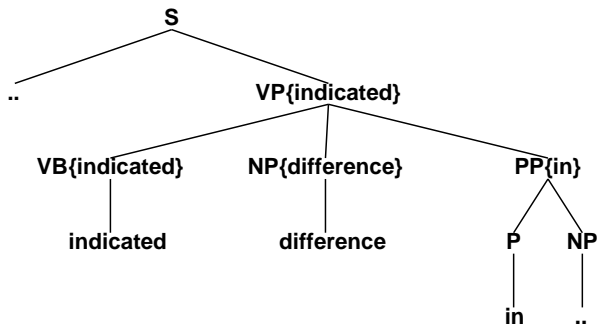
- ▶ S = a sentence

T = a parse tree

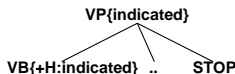
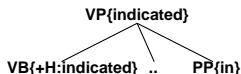
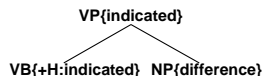
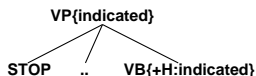
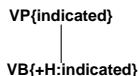
A statistical parsing model defines $P(T | S)$

- ▶ Find best parse: $\arg \max_T P(T | S)$
- ▶ $P(T | S) = \frac{P(T, S)}{P(S)} = P(T, S)$
- ▶ Best parse: $\arg \max_T P(T, S)$
- ▶ e.g. for PCFGs: $P(T, S) = \prod_{i=1 \dots n} P(\text{RHS}_i | \text{LHS}_i)$

Adding Lexical Information to PCFG



Adding Lexical Information to PCFG (Collins 99, Charniak 00)



$$\begin{aligned} &P_h(\text{VB} \mid \text{VP}, \text{indicated}) \times P_l(\text{STOP} \mid \text{VP}, \text{VB}, \text{indicated}) \times \\ &P_r(\text{NP}(\text{difference}) \mid \text{VP}, \text{VB}, \text{indicated}) \times \\ &P_r(\text{PP}(\text{in}) \mid \text{VP}, \text{VB}, \text{indicated}) \times \\ &P_r(\text{STOP} \mid \text{VP}, \text{VB}, \text{indicated}) \end{aligned}$$

Evaluation of Parsing

- ▶ Consider a candidate parse to be evaluated against the truth (or gold-standard parse):

candidate: (S (A (P this) (Q is)) (A (R a) (T test)))

gold: (S (A (P this)) (B (Q is) (A (R a) (T test))))

- ▶ In order to evaluate this, we list all the constituents

Candidate	Gold
(0,4,S)	(0,4,S)
(0,2,A)	(0,1,A)
(2,4,A)	(1,4,B)
	(2,4,A)

- ▶ Skip spans of length 1 which would be equivalent to part of speech tagging accuracy.
- ▶ Precision is defined as $\frac{\#correct}{\#proposed} = \frac{2}{3}$ and recall as $\frac{\#correct}{\#in\ gold} = \frac{2}{4}$.
- ▶ Another measure: crossing brackets,

candidate: [an [incredibly expensive] coat] (1 CB)

gold: [an [incredibly [expensive coat]]

Evaluation of Parsing

Bracketing recall R = $\frac{\text{num of correct constituents}}{\text{num of constituents in the goldfile}}$

Bracketing precision P = $\frac{\text{num of correct constituents}}{\text{num of constituents in the parsed file}}$

Complete match = % of sents where recall & precision are both 100%

Average crossing = $\frac{\text{num of constituents crossing a goldfile constituent}}{\text{num of sents}}$

No crossing = % of sents which have 0 crossing brackets

2 or less crossing = % of sents which have ≤ 2 crossing brackets

Statistical Parsing Results

System	$\leq 40wds$	$\leq 40wds$	$\leq 100wds$	$\leq 100wds$
	P	R	P	R
(Magerman 95)	84.9	84.6	84.3	84.0
(Collins 99)	88.5	88.7	88.1	88.3
(Charniak 97)	87.5	87.4	86.7	86.6
(Ratnaparkhi 97)			86.3	87.5
(Charniak 99)	90.1	90.1	89.6	89.5
(Collins 00)	90.1	90.4	89.6	89.9
Voting (HB99)	92.09	89.18		

Practical Issues: Beam Thresholding and Priors

- ▶ Probability of nonterminal X spanning $j \dots k$: $N[X, j, k]$
- ▶ Beam Thresholding compares $N[X, j, k]$ with every other Y where $N[Y, j, k]$
- ▶ But what should be compared?
- ▶ Just the *inside probability*: $P(X \overset{*}{\Rightarrow} t_j \dots t_k)$?
written as $\beta(X, j, k)$
- ▶ Perhaps $\beta(\text{FRAG}, 0, 3) > \beta(\text{NP}, 0, 3)$, but NPs are much more likely than FRAGs in general

Practical Issues: Beam Thresholding and Priors

- ▶ The correct estimate is the *outside probability*:

$$P(S \overset{*}{\Rightarrow} t_1 \dots t_{j-1} \ X \ t_{k+1} \dots t_n)$$

written as $\alpha(X, j, k)$

- ▶ Unfortunately, you can only compute $\alpha(X, j, k)$ efficiently after you finish parsing and reach $(S, 0, n)$

Practical Issues: Beam Thresholding and Priors

- ▶ To make things easier we multiply the prior probability $P(X)$ with the inside probability
- ▶ In beam Thresholding we compare every new insertion of X for span j, k as follows:
Compare $P(X) \cdot \beta(X, j, k)$ with the most probable Y
 $P(Y) \cdot \beta(Y, j, k)$
- ▶ Assume Y is the most probable entry in j, k , then we compare

$$\text{beam} \cdot P(Y) \cdot \beta(Y, j, k) \quad (1)$$

$$P(X) \cdot \beta(X, j, k) \quad (2)$$

- ▶ If $(2) < (1)$ then we prune X for this span j, k
- ▶ beam is set to a small value, say 0.001 or even 0.01.
- ▶ As the beam value increases, the parser speed increases (since more entries are pruned).
- ▶ A simpler (but not as effective) alternative to using the beam is to keep only the top K entries for each span j, k

Experiments with Beam Thresholding

