

Weighted Context-Free Grammars and Proper PCFGs

Mark Johnson

Brown University

4th August, 2005

based heavily on Zhiyi Chi (1999) “Statistical Properties of
Context-Free Grammars”, *Computational Linguistics*

Thanks to Noah Smith

Overview

- Not all PCFGs are *proper* (have a partition function $Z = 1$)
- Determining whether a PCFG is proper
- Maximum likelihood estimates of CFGs are always proper
- Weighted CFGs and Gibbs form CFGs
 - these arise naturally in conditional estimation
- WCFGs define the same distributions over trees as PCFGs
- how to convert a WCFG to an equivalent PCFG
- Conditional distributions
- CRF conditional distributions are HMM conditional distributions

Probabilistic Context-Free Grammars (1)

A PCFG is a tuple $G = (N, T, R, S, p)$ where:

- N is a finite set of nonterminals
- T is a finite set of terminals
- $R \subset N \times (N \cup T)^*$ is a finite set of rules or productions
 - assume G does not contain useless rules or symbols
- $S \in N$ is the start symbol
- p is a function from $N \times (N \cup T)^*$ to $[0, 1]$ that is non-zero only on R . $p(A \rightarrow \alpha)$ is the probability of the rule $A \rightarrow \alpha$.

$$\text{for all } A \in N, \quad \sum_{\alpha: A \rightarrow \alpha \in R} p(A \rightarrow \alpha) = 1$$

Probabilistic Context-Free Grammars (2)

The “probability” of a tree is the product of the probabilities of the rules used to generate it

$$p(t) = \prod_{A \rightarrow \alpha \in R} p(A \rightarrow \alpha)^{f_{A \rightarrow \alpha}(t)}$$

where:

- t is a parse tree and $p(t)$ is its “probability”
- $f_A(t)$ is the number of nodes labelled A in t
- $f_{A \rightarrow \alpha}(t)$ is the number of times $A \rightarrow \alpha$ is used in t

Not all PCFGs are proper

A PCFG is *proper* iff $\sum_{t \in \mathcal{T}} p(t) = 1$ where \mathcal{T} is the set of all trees

$$R = \left\{ \begin{array}{l} S \rightarrow S S \quad q \\ S \rightarrow a \quad 1 - q \end{array} \right\}$$

$$\begin{aligned} Z_h &= \sum_{t: \text{height}(t) \leq h} p(t) \\ &= (1 - q) + q s_{h-1}^2 \end{aligned}$$

so the fixed point $Z = \lim_{h \rightarrow \infty} Z_h = \sum_t p(t)$ satisfies

$$\begin{aligned} Z &= 1 - q + qZ^2 \\ Z &= \min(1, 1/q - 1) \\ &> 1 \text{ when } p > 1/2 \end{aligned}$$

Determining if a PCFG is proper

- Define a matrix M indexed by nonterminals in N

$$\begin{aligned} M_{A,B} &= \text{the expected number of } B\text{s that } A \text{ rewrites to} \\ &= \sum_{\alpha: A \rightarrow \alpha \in R} p(A \rightarrow \alpha) n_B(\alpha) \end{aligned}$$

where $n_B(\alpha)$ is the number of B s in α

- $M_{A,B}^k$ is the expected number of B s that A rewrites to in k steps
- G is proper iff each $M_{A,B}^k \rightarrow 0$ as $k \rightarrow \infty$
 - \Leftrightarrow the largest eigenvalue of M is less than 1
- Wetherell (1980) describes an efficient way of determining this

MLEs of PCFGs are always proper (1)

- Relative frequency estimator from weighted set \mathcal{D} of trees, where $W(t)$ is weight of tree t

$$p(A \rightarrow \alpha) = \frac{\sum_{t \in \mathcal{D}} f_{A \rightarrow \alpha}(t) W(t)}{\sum_{t \in \mathcal{D}} f_A(t) W(t)}$$

- encompasses MLE from treebanks, and M step of EM

$$\begin{aligned} q_A &= P(A \text{ fails to terminate}) \\ &= \sum_{A \rightarrow \alpha} p(A \rightarrow \alpha) P(\cup_i \{\alpha_i \text{ fails to terminate}\}) \\ &\leq \sum_{A \rightarrow \alpha} p(A \rightarrow \alpha) \sum_i P(\{\alpha_i \text{ fails to terminate}\}) \\ &= \sum_{A \rightarrow \alpha} p(A \rightarrow \alpha) \sum_{B \in N} n_B(\alpha) q_B \end{aligned}$$

MLEs of PCFGs are always proper (2)

$$\begin{aligned}
 q_A &\leq \sum_{A \rightarrow \alpha} p(A \rightarrow \alpha) \sum_B n_B(\alpha) q_B \\
 &= \sum_B q_B \left(\frac{\sum_{A \rightarrow \alpha} n_B(\alpha) \sum_{t \in \mathcal{D}} f_{A \rightarrow \alpha}(t) W(t)}{\sum_{t \in \mathcal{D}} f_A(t) W(t)} \right)
 \end{aligned}$$

$$\begin{aligned}
 q_A \sum_{t \in \mathcal{D}} f_A(t) W(t) &\leq \sum_B q_B \sum_{t \in \mathcal{D}} \sum_{A \rightarrow \alpha} n_B(\alpha) f_{A \rightarrow \alpha}(t) W(t) \\
 \sum_A q_A \sum_{t \in \mathcal{D}} f_A(t) W(t) &\leq \sum_B q_B \sum_{t \in \mathcal{D}} \sum_A \sum_{A \rightarrow \alpha} n_B(\alpha) f_{A \rightarrow \alpha}(t) W(t) \\
 &= \sum_B q_B \sum_{t \in \mathcal{D}} \tilde{f}_B(t) W(t)
 \end{aligned}$$

where $\tilde{f}_B(t)$ is the number of *non-root* nodes labeled B in t .

Note that $f_S(t) = \tilde{f}_S(t) + 1$, and $f_A(t) = \tilde{f}_A(t)$ for all $A \neq S$

MLEs of PCFGs are always proper (3)

$$q_A \sum_{t \in \mathcal{D}} f_A(t) W(t) \leq \sum_B q_B \sum_{t \in \mathcal{D}} \tilde{f}_B(t) W(t)$$

$$\sum_A q_A \sum_{t \in \mathcal{D}} \left(f_A(t) - \tilde{f}_A(t) \right) W(t) \leq 0$$

But since $f_A(t) = \tilde{f}_A(t)$ for $A \neq S$ and $f_S(t) = \tilde{f}_S(t) + 1$:

$$q_S \sum_t W(t) \leq 0$$

which implies that $q_S = 0$

Weighted CFGs

- A *weighted CFG* G is one where each rule $A \rightarrow \alpha \in R$ is associated with a positive weight $w_{A \rightarrow \alpha}$. The weight $w(t)$ of a tree t generated by G is the product of the weights of the rules that generate it.

$$w(t) = \prod_{A \rightarrow \alpha \in R} w_{A \rightarrow \alpha}^{f_{A \rightarrow \alpha}(t)}$$

$$Z = \sum_{t \in \mathcal{T}} w(t)$$

$$P(t) = w(t)/Z$$

- WCFGs can also be expressed as *Gibbs* or *log linear* models

$$\lambda_{A \rightarrow \alpha} = \log w_{A \rightarrow \alpha} \text{ for each } A \rightarrow \alpha \in R$$

$$w(t) = \prod_{A \rightarrow \alpha} \exp(\lambda_{A \rightarrow \alpha})^{f_{A \rightarrow \alpha}(t)} = \exp \sum_{A \rightarrow \alpha} \lambda_{A \rightarrow \alpha} f_{A \rightarrow \alpha}(t)$$

$$P(t) = \frac{1}{Z} \exp \sum_{A \rightarrow \alpha} \lambda_{A \rightarrow \alpha} f_{A \rightarrow \alpha}(t)$$

Reasons for using WCFGs

- Unconstrained numerical optimization is easier than constrained numerical optimization
 - ⇒ numerically estimating a WCFG is easier than a PCFG
- Conditional estimation (trees given strings)
 - should be more accurate for parsing (given our bad grammars)
 - no closed form known (to me) ⇒ numerical methods
 - there is a dynamic programming algorithm for calculating conditional likelihood and its derivatives
- Want to impose a prior on rule probabilities or regularize
 - Except for Dirichlet prior (which is conjugate to multinomial), numerical optimization probably required
 - Often easier to state priors/regularizers in log linear parameter space $\lambda_{A \rightarrow \alpha}$

Every WCFG dist is a proper PCFG dist

- In terms of grammars, WCFGs \supset PCFGs \supset proper PCFGs, but these all define exactly the same probability distributions over trees
- Chi 1999's rule probabilities for the equivalent proper PCFG:

$$p(A \rightarrow \alpha) = \frac{1}{Z_A} w_{A \rightarrow \alpha} \prod_{k=1}^{|\alpha|} Z_{\alpha_k}, \text{ where}$$

$$Z_A = \sum_{t \in \mathcal{T}_A} w(t)$$

$$\mathcal{T}_A = \text{set of trees rooted in } A$$

$$P_B(t) = \frac{1}{Z_B} \prod_{A \rightarrow \alpha \in R} w_{A \rightarrow \alpha}^{f_{A \rightarrow \alpha}(t)} \text{ for } t \in \mathcal{T}_B$$

\Rightarrow many different WCFGs define the same distribution over trees

\Rightarrow WCFGs weights are not identifiable even from treebank data

WCFG to PCFG conversion (1)

Proposition: $p_B(t) = P_B(t)$ for all $B \in N$ and $t \in \mathcal{T}_B$, where

$$\begin{aligned} p(A \rightarrow \alpha) &= \frac{1}{Z_A} w_{A \rightarrow \alpha} \prod_{k=1}^{|\alpha|} Z_{\alpha_k} \\ p_B(t) &= \prod_{A \rightarrow \alpha} p(A \rightarrow \alpha)^{f_{A \rightarrow \alpha}(t)} \\ P_B(t) &= \frac{1}{Z_B} \prod_{A \rightarrow \alpha} w_{A \rightarrow \alpha}^{f_{A \rightarrow \alpha}(t)} \end{aligned}$$

Proof by induction on the height h of t . Trivial for $h = 1$ (terminals)

WCFG to PCFG conversion (2)

Suppose $t \in \mathcal{T}_B$ has height $h > 1$, root rule is $B \rightarrow \beta$ and t_k is subtree rooted in k th child of t .

$$\begin{aligned} P_B(t) &= \frac{1}{Z_B} \prod_{A \rightarrow \alpha} w_{A \rightarrow \alpha} f_{A \rightarrow \alpha}(t) \\ &= \frac{1}{Z_B} w_{B \rightarrow \beta} \prod_{k=1}^{|\beta|} \prod_{A \rightarrow \alpha} w_{A \rightarrow \alpha} f_{A \rightarrow \alpha}(t_k) \\ &= \frac{1}{Z_B} w_{B \rightarrow \beta} \prod_{k=1}^{|\beta|} Z_{\beta_k} P_{\beta_k}(t_k) \\ &= \frac{1}{Z_B} w_{B \rightarrow \beta} \prod_{k=1}^{|\beta|} Z_{\beta_k} p_{\beta_k}(t_k) \text{ by induction hyp} \\ &= p(B \rightarrow \beta) \prod_{k=1}^{|\beta|} p_{\beta_k}(t_k) \\ &= p_B(t) \end{aligned}$$

Calculating WCFG partition functions Z_A

$$p(A \rightarrow \alpha) = \frac{1}{Z_A} w_{A \rightarrow \alpha} \prod_{k=1}^{|\alpha|} Z_{\alpha_k}$$

Let $\mathcal{T}_A^{(h)}$ = trees rooted in A of height $\leq h$

$$Z_A^{(h)} = \sum_{t \in \mathcal{T}_A^{(h)}} \prod_{A \rightarrow \alpha} w_{A \rightarrow \alpha}^{f_{A \rightarrow \alpha}(t)}$$

$$Z_A = \lim_{h \rightarrow \infty} Z_A^{(h)}$$

$$Z_A^{(1)} = 1 \text{ if } A \in T, 0 \text{ otherwise}$$

$$Z_A^{(h+1)} = \begin{cases} 1 & \text{for } A \in T \\ \sum_{\alpha: A \rightarrow \alpha \in R} w_{A \rightarrow \alpha} \prod_{k=1}^{|\alpha|} Z_{\alpha_k}^{(h)} & \text{for } A \in N \end{cases}$$

Conditional Distributions

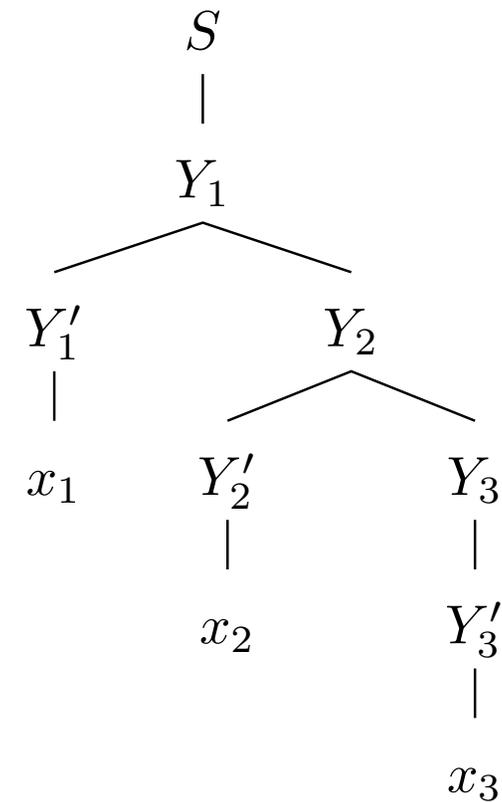
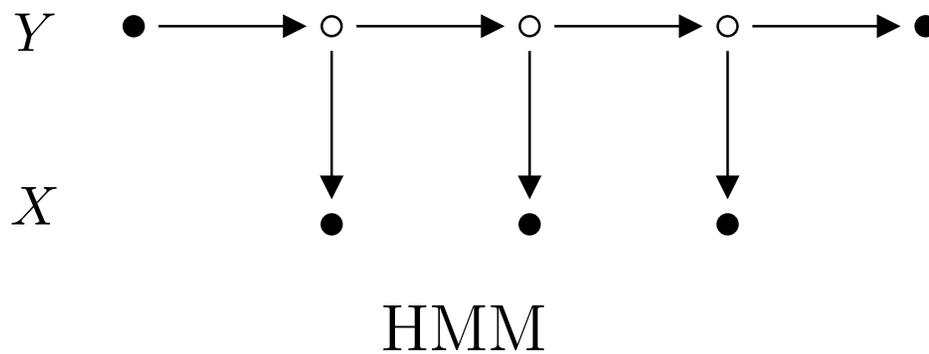
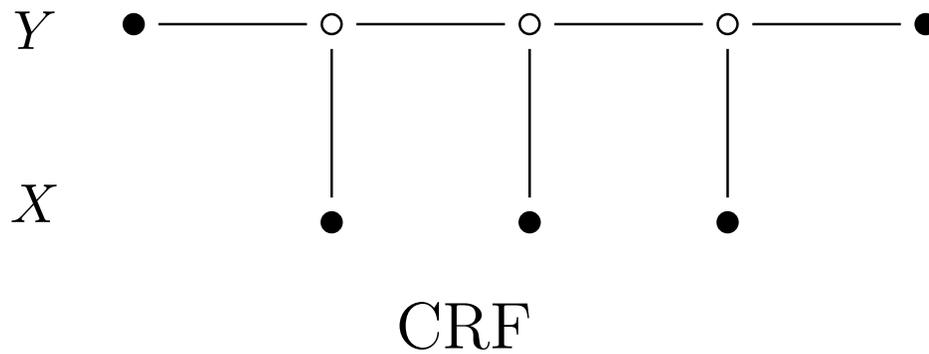
- Some WCFGs define *conditional* distributions of trees given strings, even though their partition functions Z diverge

$$S \rightarrow S S; 1 \quad S \rightarrow a; 1$$

- Chi's formula requires that all partition functions Z converge
- Change the WCFG weights $w_{A \rightarrow \alpha}$ to $\gamma^{|\alpha|-1} w_{A \rightarrow \alpha}$, $\gamma > 0$
 - multiplies the weight of a tree with yield y by $\gamma^{|y|-1}$ \Rightarrow doesn't affect the conditional distribution
- To find a PCFG with same conditional distribution as a WCFG
 1. Search for γ that makes the WCFG partition functions converge
 2. Apply the Chi formula to obtain PCFG rule probabilities

CRF cond dists = HMM cond dists

- A CRF is a WCFG $\Rightarrow \exists$ HMM with same cond dist



Conclusion

- Not all PCFGs are *proper* (have a partition function $Z = 1$)
- Determining whether a PCFG is proper
- Maximum likelihood estimates of CFGs are always proper
- Weighted CFGs and Gibbs form CFGs
 - these arise naturally in conditional estimation
- WCFGs define the same distributions over trees as PCFGs
- how to convert a WCFG to an equivalent PCFG
- Conditional distributions
- CRF conditional distributions are HMM conditional distributions