

Count-Based State Merging for Probabilistic Regular Tree Grammars

Toni Dietze

Faculty of Computer Science
Technische Universität Dresden
01062 Dresden, Germany
`toni.dietze@tu-dresden.de`

Mark-Jan Nederhof

School of Computer Science
University of St Andrews
KY16 9SX, UK

Abstract

We present an approach to obtain language models from a tree corpus using probabilistic regular tree grammars (prt_g). Starting with a prt_g only generating trees from the corpus, the prt_g is generalized step by step by merging nonterminals. We focus on bottom-up deterministic prt_g to simplify the calculations.

1 Introduction

Constituent parsing plays an important role in natural language processing (nlp). One can easily read off a pcfg from a tree corpus and use it for parsing. This might work quite well (Charniak, 1996), but it can be even more fruitful to introduce a state behaviour that is not visible in the corpus (Klein and Manning, 2003). The Expectation-Maximization Algorithm (Dempster et al., 1977) can be used to train probabilities if the state behaviour is fixed (Matsuzaki et al., 2005). This can be improved by adapting the state behaviour automatically by cleverly splitting and merging states (Petrov et al., 2006).

More generally, finding a grammar by examining terminal objects is one of the problems investigated in the field of grammatical inference. There are many results for the string case, e.g., on how to learn deterministic stochastic finite (string) automata from text (Carrasco and Oncina, 1994; Carrasco and Oncina, 1999). For the tree case, there are, e.g., results for identifying function distinguishable regular tree languages from text (Fernau, 2002). There is also a generalization of n -grams to trees including smoothing techniques (Rico-Juan et al., 2000; Rico-Juan et al., 2002).

The mentioned results for deterministic stochastic finite (string) automata were generalized to an algorithm that learns stochastic deterministic tree automata from trees (Carrasco et al., 2001). Given

a tree corpus, this approach yields a single grammar. The authors experimentally showed that if the corpus is too small, the grammar tends to be too general. In contrast, the split-merge approach of Petrov et al. (2006) produces a sequence of different grammars. One can use cross validation to select a grammar from that sequence that suitably abstracts away from the training corpus. Because of the intricate combination of splitting and merging however, the behavior is very difficult to analyse theoretically. Our approach is similar to the split-merge procedure in that a sequence of grammars is obtained. However, it differs by operating without splitting and relying on merging alone, thereby obtaining a simpler framework.

Our goal is to create a sequence of probabilistic regular tree grammars (prt_g) from a corpus such that every prt_g abstracts away from the corpus more than its predecessors in the sequence (cf. Algorithm 1). We start with a prt_g that admits no more and no less than the trees in the corpus. The rules of the prt_g are then changed step by step to make the grammar more general.

We generalize a prt_g by *merging* nonterminals, which means we replace several nonterminals by a single new one. The weights of the resulting prt_g are assigned by maximum likelihood estimation on the corpus.

To make the approach easier, we only consider bottom-up deterministic prt_g. On the one hand, this simplifies our calculations, e.g., the maximum likelihood estimate; on the other hand, this simplifies the application of the prt_g as language model, since the search for the most probable tree for a given yield does not have to consider several derivations for a single tree.

2 Preliminaries

Let X be a set. The *identity relation on X* is defined as $\{(x, x) \mid x \in X\}$. An *equivalence relation on X* is a reflexive, symmetric, and transitive

relation $(\equiv) \subseteq X \times X$. We write $x \equiv y$ instead of $(x, y) \in (\equiv)$ for $x, y \in X$. Note that the identity relation is an equivalence relation. Let $x \in X$. The *equivalence class of x (induced by (\equiv))*, denoted by $[x]_{\equiv}$ (or just $[x]$ if (\equiv) is clear from the context), is defined as $\{y \in X \mid x \equiv y\}$. The *quotient set of X by (\equiv)* , denoted by X/\equiv , is defined as $\{[x] \mid x \in X\}$. Note that X/\equiv is a partition of X and that for every partition P of X there is an equivalence relation (\equiv') such that $P = X/\equiv'$. The set of equivalence relations over X forms a complete lattice ordered by set inclusion.

We denote the set of the natural numbers including 0 by \mathbb{N} . An *alphabet* (of *symbols*) is a finite, non-empty set. A *ranked alphabet* is an alphabet Σ where we associate a *rank* $\text{rk}(\sigma) \in \mathbb{N}$ with every $\sigma \in \Sigma$. Let Σ be a ranked alphabet. The *set of trees over Σ* , denoted by \mathbb{T}_{Σ} , is the smallest set T such that $\sigma(t_1, \dots, t_{\text{rk}(\sigma)}) \in T$ for every $\sigma \in \Sigma$ and $t_1, \dots, t_{\text{rk}(\sigma)} \in T$. Additionally, we define $\mathbb{T}_{\emptyset} = \emptyset$. Let $t = \sigma(t_1, \dots, t_k) \in \mathbb{T}_{\Sigma}$. The *set of positions of t* , denoted by $\text{pos}(t)$, is defined as $\{\varepsilon\} \cup \{iw \mid i \in \{1, \dots, k\}, w \in \text{pos}(t_i)\}$. Let $w \in \text{pos}(t)$. The *symbol of t at w* , denoted by $t(w)$, is defined as σ if $w = \varepsilon$, and by $t_i(w')$ if $w = iw'$. The *subtree of t at w* , denoted by $t|_w$, is defined as t if $w = \varepsilon$, and by $t_i|_{w'}$ if $w = iw'$.

Let X be a set and $f: X \rightarrow \mathbb{R}_{\geq 0}$ a mapping. The *support of f* , denoted by $\text{supp}(f)$, is defined as $\{x \in X \mid f(x) \neq 0\}$. The *size of f* , denoted by $|f|$, is defined as $\sum_{x \in \text{supp}(f)} f(x)$. We call f a *corpus* (over X) if $\text{supp}(f)$ is finite and non-empty. We call f a *probability distribution* (over X) if $|f| = 1$. We denote the set of all probability distributions over X by $\text{Pd}(X)$. Note that f may be both a corpus and a probability distribution. Sometimes we refer to the values of a corpus by the word *counts*.

Definition 1. A *regular tree grammar* (rtg) is a tuple (N, Σ, I, R) where

- N is an alphabet (of *nonterminals*),
- Σ is a ranked alphabet (of *terminals*) such that $N \cap \Sigma = \emptyset$,
- $I \subseteq N$ is a non-empty set (of *initial nonterminals*), and
- R is a finite set of *rules* of the form $A_0 \rightarrow \sigma(A_1, \dots, A_{\text{rk}(\sigma)})$ where $\sigma \in \Sigma$ and $A_0, \dots, A_{\text{rk}(\sigma)} \in N$. \square

Our definition above corresponds to a normal form of rtg with regard to a more general definition, which was given by, for example, Gécseg and

Steinby (1984, Chapter II, Section 3).

Let $G = (N, \Sigma, I, R)$ be an rtg. For a rule $r \in R$ of the form $A_0 \rightarrow \sigma(A_1, \dots, A_{\text{rk}(\sigma)})$ we define $\text{lhs}(r) = A_0$ and $\text{rhs}(r) = \sigma(A_1, \dots, A_{\text{rk}(\sigma)})$. We associate with r the rank $\text{rk}(r) = \text{rk}(\sigma)$, hence we may view R as a ranked alphabet if R is non-empty. We call G *bottom-up deterministic* if $\text{rhs}(r_1) = \text{rhs}(r_2)$ implies $r_1 = r_2$ for every $r_1, r_2 \in R$.

Let $t \in \mathbb{T}_{\Sigma}$. A *derivation tree of G for t* is a tree $d \in \mathbb{T}_R$ such that $\text{pos}(d) = \text{pos}(t)$, $\text{lhs}(d(\varepsilon)) \in I$, and for every $w \in \text{pos}(t)$ we have $\text{rhs}(d(w)) = t(w)(\text{lhs}(d(w_1)), \dots, \text{lhs}(d(w_{\text{rk}(t(w))})))$. The *set of all derivation trees of G for t* is denoted by $D_G(t)$. The *language of trees generated by G* , denoted by $\llbracket G \rrbracket$, is defined as $\{t \in \mathbb{T}_{\Sigma} \mid D_G(t) \neq \emptyset\}$. Note that if G is bottom-up deterministic, $D_G(t)$ has at most one element. We denote this element, if it exists, by d_G^t .

Definition 2. A *probabilistic regular tree grammar* (prtg) is a tuple (G, ι, ρ) where

- $G = (N, \Sigma, I, R)$ is an rtg,
- $\iota: I \rightarrow [0, 1]$ is a mapping (*initial weights*), and
- $\rho: R \rightarrow [0, 1]$ is a mapping (*rule weights*). \square

Let $P = (G, \iota, \rho)$ be a prtg. Terminology for rtg is carried over to prtg, e.g., P is bottom-up deterministic iff G is bottom-up deterministic. We call P *proper* if ι is a probability distribution and $\sum_{r \in R: \text{lhs}(r)=A} \rho(r) = 1$ for every $A \in N$.

Let $t \in \mathbb{T}_{\Sigma}$ and $d \in D_P(t)$. The *weight of d (induced by P)*, denoted by $\llbracket P \rrbracket(d)$, is defined as $\iota(\text{lhs}(d(\varepsilon))) \cdot \prod_{w \in \text{pos}(t)} \rho(d(w))$. The *weight of t (induced by P)*, denoted by $\llbracket P \rrbracket(t)$, is defined as $\sum_{d \in D_P(t)} \llbracket P \rrbracket(d)$. If $\llbracket P \rrbracket$ is a probability distribution over \mathbb{T}_{Σ} , then P is called *consistent*.

Let c be a corpus over \mathbb{T}_{Σ} . The *canonical rtg of c* and the *canonical prtg of c* are defined as $G = (N, \Sigma, I, R)$ and $P = (G, \iota, \rho)$, respectively, where

- $N = \{t|_w \mid t \in \text{supp}(c), w \in \text{pos}(t)\}$,
- $I = \text{supp}(c)$,
- $R = \{t|_{\varepsilon} \rightarrow t(\varepsilon)(t|_1, \dots, t|_{\text{rk}(t(\varepsilon))}) \mid t \in N\}$,
- $\iota(t) = \frac{c(t)}{|c|}$ for every $t \in I$, and
- $\rho(r) = 1$ for every $r \in R$.

Note that every canonical prtg is bottom-up deterministic, proper, and consistent, and that $\llbracket P \rrbracket(t) = \iota(t)$ for every $t \in \text{supp}(\llbracket P \rrbracket)$; hence, $\llbracket P \rrbracket$ represents the relative frequencies of the trees in c .

Let $c: X \rightarrow \mathbb{R}_{\geq 0}$ be a corpus and $p \in \text{Pd}(X)$.

The *likelihood* of c given p is defined as

$$L(c | p) = \prod_{t \in \text{supp}(c)} p(t)^{c(t)}.$$

Let $\mathcal{M} \subseteq \text{Pd}(X)$. The *maximum likelihood estimate* from \mathcal{M} for c , denoted by $\text{mle}_{\mathcal{M}}(c)$, is defined as

$$\text{mle}_{\mathcal{M}}(c) = \underset{p \in \mathcal{M}}{\text{argmax}} L(c | p).$$

If $\mathcal{M} = \text{Pd}(\text{T}_{\Sigma})$, then $\text{mle}_{\mathcal{M}}(c)$ maps a tree to its relative frequency in c (Prescher, 2004, Theorem 1). Hence, the canonical prt \bar{g} represents the corpus perfectly.

Let $G = (N, \Sigma, I, R)$ be an rt \bar{g} , $\mathcal{M} = \{\llbracket P \rrbracket | P = (G, \iota, \rho) \text{ is a consistent prt}\bar{g}\}$, and c a corpus over T_{Σ} . Then we also write $\text{mle}_G(c)$ instead of $\text{mle}_{\mathcal{M}}(c)$.

Let $G = (N, \Sigma, I, R)$ be a bottom-up deterministic rt \bar{g} and c a corpus over T_{Σ} such that $\text{supp}(c) \subseteq \llbracket G \rrbracket$. Note that there is exactly one derivation tree d_G^t of G for every tree $t \in \text{supp}(c)$. Based on G , we derive three corpora from c :

$$c_G^R: R \rightarrow \mathbb{R}_{\geq 0}: r \mapsto \sum_{t \in \text{supp}(c)} c(t) \cdot |\{w \in \text{pos}(t) | r = d_G^t(w)\}|,$$

$$c_G^N: N \rightarrow \mathbb{R}_{\geq 0}: A \mapsto \sum_{t \in \text{supp}(c)} c(t) \cdot |\{w \in \text{pos}(t) | A = \text{lhs}(d_G^t(w))\}|,$$

$$c_G^I: N \rightarrow \mathbb{R}_{\geq 0}: A \mapsto \sum_{\substack{t \in \text{supp}(c): \\ A = \text{lhs}(d_G^t(\varepsilon))}} c(t).$$

Now $\text{mle}_G(c) = \llbracket (G, \iota, \rho) \rrbracket$ where for every $A \in I$ and $r \in R$ (Prescher, 2004, Theorem 10):

$$\iota(A) = \frac{c_G^I(A)}{|c|} \quad \text{and} \quad \rho(r) = \frac{c_G^R(r)}{c_G^N(\text{lhs}(r))}.$$

Note that $c_G^N(A) = \sum_{r \in R: A = \text{lhs}(r)} c_G^R(r)$.

3 Count-Based State Merging

Algorithm 1 summarizes the idea of our approach. We detail the used notions in what follows.

Let $G = (N, \Sigma, I, R)$ be an rt \bar{g} and (\equiv) an equivalence relation on N . The G -merger w.r.t. (\equiv) is the overloaded expression π_{\equiv} for

- nonterminals: $\pi_{\equiv}(A) = [A]$ for every $A \in N$,
- rules: $\forall \sigma \in \Sigma: \forall A_0, \dots, A_{\text{rk}(\sigma)} \in N:$
 $\pi_{\equiv}(A_0 \rightarrow \sigma(A_1, \dots, A_{\text{rk}(\sigma)}))$

Algorithm 1 Count-Based State Merging

Input: corpus c over T_{Σ}

Output: sequence of bottom-up deterministic prt \bar{g} P_0, \dots, P_n , some $n \in \mathbb{N}$, such that $\text{supp}(\llbracket P_0 \rrbracket) \subseteq \dots \subseteq \text{supp}(\llbracket P_n \rrbracket)$

- 1: $P_0 = (G_0, \iota_0, \rho_0) \leftarrow$ canonical prt \bar{g} of c
 - 2: $i \leftarrow 0$
 - 3: **while** there exists a non-trivial G_i -merger **do**
 - 4: $\pi \leftarrow \text{BESTMERGER}(G_i, c)$
 - 5: $i \leftarrow i + 1$
 - 6: $G_i \leftarrow \pi(G_{i-1})$
 - 7: let P_i be prt \bar{g} such that $\text{mle}_{G_i}(c) = \llbracket P_i \rrbracket$
-

$$= [A_0] \rightarrow \sigma([A_1], \dots, [A_{\text{rk}(\sigma)}]),$$

- sets of nonterminals or rules by applying π_{\equiv} elementwise, and

- rt \bar{g} : $\pi_{\equiv}(G) = (N/\equiv, \Sigma, \pi_{\equiv}(I), \pi_{\equiv}(R))$.

Note that $\llbracket G \rrbracket \subseteq \llbracket \pi_{\equiv}(G) \rrbracket$, because by replacing each rule r in a derivation tree of G by $\pi_{\equiv}(r)$ we get a derivation tree of $\pi_{\equiv}(G)$. We call π_{\equiv} *non-trivial*, if (\equiv) is not the identity relation. We say π_{\equiv} *merges* A_1 and A_2 , iff $A_1 \equiv A_2$. We carry over the lattice structure of the set of equivalence relations over N to the set of G -mergers in order to identify minimal and least G -mergers with certain properties.

To deal with prt \bar{g} , we fix a corpus c over T_{Σ} such that $\text{supp}(c) \subseteq \llbracket G \rrbracket$. We repeatedly use the maximum likelihood estimate to assign weights to an rt \bar{g} . That is, the weights are not manipulated during merging itself, but they are used to choose a G -merger: Let Π be a set of G -mergers. The *best G -merger from Π w.r.t. c* is defined as

$$\underset{\pi \in \Pi}{\text{argmax}} L(c | \llbracket \text{mle}_{\pi(G)}(c) \rrbracket). \quad (1)$$

So far, the presented notions are defined for general rt \bar{g} . Now let $G = (N, \Sigma, I, R)$ be a bottom-up deterministic rt \bar{g} . Assuming $0^0 = 1$, we then have the following, which is proven in Appendix A:

$$L(c | \llbracket \text{mle}_G(c) \rrbracket) = \frac{\prod_{A \in N} c_G^I(A) c_G^I(A)}{|c|^{|c|}} \cdot \frac{\prod_{r \in R} c_G^R(r) c_G^R(r)}{\prod_{A \in N} c_G^N(A) c_G^N(A)}. \quad (2)$$

This can be used to make Expression 1 more manageable: Let $G = (N, \Sigma, I, R)$ be a bottom-up deterministic rt \bar{g} and let Π be the set of all non-trivial G -mergers π such that $\pi(G)$ is bottom-up deterministic. Then Equation 2 gives a more direct way of computing the likelihood in Express-

sion 1 without explicitly calculating $\text{mle}_{\pi(G)}(c)$. This is the reason for calling our approach “count-based”, and this is the idea of BESTMERGER in Algorithm 1.

Improving efficiency Let $\pi \in \Pi$ and $G' = (N', \Sigma, I', R') = \pi(G)$. We need $c_{G'}^R$, $c_{G'}^N$, and $c_{G'}^I$ to calculate $L(c \mid \llbracket \text{mle}_{G'}(c) \rrbracket)$. Bottom-up determinism allows us to derive $d_{G'}^t$ for every $t \in \text{supp}(c)$ by replacing every rule r in d_G^t by $\pi(r)$. Hence,

$$\forall x' \in X': c_{G'}^X(x') = \sum_{x \in X: x' = \pi(x)} c_G^X(x), \quad (3)$$

where X is any of R , N , or I (with or without prime, italic or roman). So we can reuse the corpora related to G to calculate $L(c \mid \llbracket \text{mle}_{G'}(c) \rrbracket)$.

We may rewrite Expression 1 by dividing the likelihood by $L(c \mid \llbracket \text{mle}_G(c) \rrbracket)$. Then, for many instantiations of π many factors in the fraction cancel out. In detail: Let (\equiv) be the equivalence relation underlying π , and

$$\begin{aligned} \overline{N} &= \{A \in N \mid |[A]| > 1\}, & \overline{N}' &= \pi(\overline{N}), \\ \overline{R} &= \{r \in R \mid |[r]| > 1\}, & \overline{R}' &= \pi(\overline{R}), \end{aligned}$$

where (\equiv) is extended to rules such that $r_1 \equiv r_2$ iff $\pi(r_1) = \pi(r_2)$ for every $r_1, r_2 \in R$. Then, with Equation 2 and $G' = \pi(G)$:

$$\begin{aligned} \frac{L(c \mid \llbracket \text{mle}_{G'}(c) \rrbracket)}{L(c \mid \llbracket \text{mle}_G(c) \rrbracket)} &= \frac{\prod_{A \in \overline{N}'} c_{G'}^I(A) c_{G'}^I(A)}{\prod_{A \in \overline{N}} c_G^I(A) c_G^I(A)} \\ &\cdot \frac{\prod_{r \in \overline{R}'} c_{G'}^R(r) c_{G'}^R(r)}{\prod_{r \in \overline{R}} c_G^R(r) c_G^R(r)} \cdot \frac{\prod_{A \in \overline{N}} c_G^N(A) c_G^N(A)}{\prod_{A \in \overline{N}'} c_{G'}^N(A) c_{G'}^N(A)}. \end{aligned} \quad (4)$$

Yet, finding the maximum is still expensive.

Heuristics Assume $\overline{N} = \{A_1, A_2\}$, $\overline{R} = \emptyset$ and $\overline{N} \cap I = \emptyset$. Then, using Equation 3, the right-hand side of Equation 4 is equal to $f(c_G^N(A_1), c_G^N(A_2))$ where

$$f(x, y) = \frac{x^x \cdot y^y}{(x + y)^{x+y}}. \quad (5)$$

For positive x and y , $f(x, y)$ is monotonically decreasing (cf. Appendix B). Hence, with our assumption, it is best to merge nonterminals with the least counts w.r.t. c_G^N .

This may be used to guide the search for the best merger. Recall that we want to generalize the canonical (p)rtg step by step. We want to take the smallest steps possible w.r.t. loss of likelihood

Algorithm 2

Input: rtg $G = (N, \Sigma, I, R)$ and equivalence relation (\equiv_0) on N

Output: the least G -merger π_{\equiv} such that $\pi_{\equiv}(G)$ is bottom-up deterministic and $(\equiv_0) \subseteq (\equiv)$

- 1: $(\equiv) \leftarrow (\equiv_0)$
 - 2: **while** $\pi_{\equiv}(G)$ not bottom-up deterministic **do**
 - 3: find rules r_1 and r_2 in G such that
 $\text{rhs}(\pi_{\equiv}(r_1)) = \text{rhs}(\pi_{\equiv}(r_2))$, but
 $\text{lhs}(\pi_{\equiv}(r_1)) \neq \text{lhs}(\pi_{\equiv}(r_2))$
 - 4: let $A_1 = \text{lhs}(r_1)$ and $A_2 = \text{lhs}(r_2)$
 - 5: let (\equiv') equivalence relation s.t. $N/\equiv' = N/\equiv \setminus \{[A_1], [A_2]\} \cup \{[A_1] \cup [A_2]\}$
 - 6: replace (\equiv) by (\equiv')
-

(cf. Expression 1), so we consider only minimal non-trivial mergers, i.e., mergers that merge exactly two nonterminals. Note that the application of larger mergers may be decomposed into a sequential application of several minimal non-trivial mergers. We can easily sort the minimal non-trivial mergers by the value of f for the counts of the merged nonterminals. Note that the merger which merges the nonterminals with the lowest counts comes first. We choose a beam width $n > 0$ and select the n first mergers for further investigation assuming that the best merger is among them.

Unfortunately, a minimal non-trivial merger π_{\equiv} does not necessarily result in a bottom-up deterministic rtg, but there is a least (i.e. unique minimal) merger $\pi_{\equiv'}$ such that $(\equiv) \subseteq (\equiv')$ which does (cf. Lemma 2 in Appendix C). We use Algorithm 2 to calculate this merger for each of the n chosen mergers and determine the best merger from the results w.r.t. Equation 4.

To restrict the number of considered mergers even more, it may be useful to only merge nonterminals which produce the same terminal. Note that Algorithm 2 preserves this property.

4 Practical results and outlook

So far, our implementation is not competitive. We are only able to train with small corpora. For example, training with 5 900 trees (consisting of 120 000 distinct subtrees) and a beam width of 1 000 takes about 8 days. For some inputs Algorithm 2 is very expensive. It remains to be seen whether we can avoid such inputs or reduce the effort by reusing results from previous iterations in Algorithm 1.

There is generally a very large number of non-terminals with count 1 in the canonical rtg. Our current heuristics yields the same value for every merger which considers two of such nonterminals. Especially in the first iterations of Algorithm 1 the number of mergers with the same (lowest) heuristic value far exceeds the beam width. This means it is arbitrary which mergers lie within the beam. We hope to improve the heuristics by comparing the trees which are produced by the merged non-terminals.

Using the generated grammars for parsing, we are currently only able to process sentences consisting of words seen in the training data. Even for this limited subset of sentences we are not able to improve precision or recall of brackets (Sekine and Collins, 1997) in comparison to the probabilistic context-free grammar straightforwardly obtained from the corpus. We hope this will improve with a better heuristics.

We have restricted our attention to bottom-up deterministic regular tree grammars. Thanks to this, the conceptual framework could remain relatively simple. What is unclear at this time is whether the bottom-up determinism *per se* restricts the potential accuracy of the models, in relation to the split-merge framework, which allows nondeterministic regular tree grammars.

References

- Rafael C. Carrasco and Jose Oncina. 1994. Learning stochastic regular grammars by means of a state merging method. In Rafael C. Carrasco and Jose Oncina, editors, *Grammatical Inference and Applications*, volume 862 of *Lecture Notes in Computer Science*, page 139–152. Springer Berlin Heidelberg.
- Rafael C. Carrasco and Jose Oncina. 1999. Learning deterministic regular grammars from stochastic samples in polynomial time. *RAIRO – Theoretical Informatics and Applications*, 33(1):1–19.
- Rafael C. Carrasco, Jose Oncina, and Jorge Calera-Rubio. 2001. Stochastic inference of regular tree languages. *Machine Learning*, 44(1-2):185–197.
- Eugene Charniak. 1996. Tree-bank grammars. In *Proc. of AAAI/IAAI 1996*, volume 2, pages 1031–1036.
- A. P. Dempster, N. M. Laird, and D. B. Rubin. 1977. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, 39(1):1–38.
- Henning Fernau. 2002. Learning tree languages from text. In Jyrki Kivinen and Robert H. Sloan, editors, *Computational Learning Theory*, volume 2375 of *Lecture Notes in Computer Science*, page 153–168. Springer Berlin Heidelberg.
- Ferenc Gécseg and Magnus Steinby. 1984. *Tree Automata*. Akadémiai Kiadó, Budapest, Hungary.
- Dan Klein and Christopher D. Manning. 2003. Accurate unlexicalized parsing. In *Proc. of ACL 2003*, volume 1, pages 423–430.
- Takuya Matsuzaki, Yusuke Miyao, and Jun’ichi Tsujii. 2005. Probabilistic CFG with latent annotations. In *Proc. of ACL 2005*, pages 75–82.
- Slav Petrov, Leon Barrett, Romain Thibaux, and Dan Klein. 2006. Learning accurate, compact, and interpretable tree annotation. In *Proc. of COLING/ACL 2006*, pages 433–440.
- Detlef Prescher. 2004. A tutorial on the expectation-maximization algorithm including maximum-likelihood estimation and EM training of probabilistic context-free grammars. *CoRR*, abs/cs/0412015.
- Juan Ramón Rico-Juan, Jorge Calera-Rubio, and Rafael C. Carrasco. 2000. Probabilistic k -testable tree languages. In Arlindo L. Oliveira, editor, *Grammatical Inference: Algorithms and Applications*, volume 1891 of *Lecture Notes in Computer Science*, page 221–228. Springer Berlin Heidelberg.
- Juan Ramón Rico-Juan, Jorge Calera-Rubio, and Rafael C. Carrasco. 2002. Stochastic k -testable tree languages and applications. In Pieter Adriaans, Henning Fernau, and Menno van Zaanen, editors, *Grammatical Inference: Algorithms and Applications*, volume 2484 of *Lecture Notes in Computer Science*, page 199–212. Springer Berlin Heidelberg.
- Satoshi Sekine and Michael John Collins. 1997. Evalb. <http://nlp.cs.nyu.edu/evalb/>. Accessed 2015-03-19.

A The likelihood of the maximum likelihood estimate in the bottom-up deterministic case (Equation 2)

Let $G = (N, \Sigma, I, R)$ be a bottom-up deterministic rtg and c a corpus over T_Σ such that $\text{supp}(c) \subseteq \llbracket G \rrbracket$. Let $P = (G, \iota, \rho)$ such that $\llbracket P \rrbracket = \text{mle}_G(c)$. We can transform $L(c \mid \llbracket P \rrbracket)$ as follows, assuming $0^0 = 1$:

$$\begin{aligned}
& L(c \mid \llbracket P \rrbracket) \\
&= \prod_{t \in \text{supp}(c)} \llbracket P \rrbracket(t)^{c(t)} && \text{(def. of } L) \\
&= \prod_{t \in \text{supp}(c)} \left(\sum_{d \in D_P(t)} \llbracket P \rrbracket(d) \right)^{c(t)} && \text{(def. of } \llbracket P \rrbracket) \\
&= \prod_{t \in \text{supp}(c)} \llbracket P \rrbracket(d_G^t)^{c(t)} && \text{(assumptions for } G \text{ and } c) \\
&= \prod_{t \in \text{supp}(c)} \left(\iota(\text{lhs}(d_G^t(\varepsilon))) \cdot \prod_{w \in \text{pos}(t)} \rho(d_G^t(w)) \right)^{c(t)} && \text{(def. of } \llbracket P \rrbracket) \\
&= \prod_{t \in \text{supp}(c)} \iota(\text{lhs}(d_G^t(\varepsilon)))^{c(t)} \cdot \prod_{w \in \text{pos}(t)} \rho(d_G^t(w))^{c(t)} && \text{(distributivity)} \\
&= \left(\prod_{t \in \text{supp}(c)} \iota(\text{lhs}(d_G^t(\varepsilon)))^{c(t)} \right) \cdot \prod_{t \in \text{supp}(c)} \prod_{w \in \text{pos}(t)} \rho(d_G^t(w))^{c(t)} && \text{(commutativity)} \\
&= \left(\prod_{A \in N} \prod_{\substack{t \in \text{supp}(c): \\ A = \text{lhs}(d_G^t(\varepsilon))}} \iota(A)^{c(t)} \right) \cdot \prod_{r \in R} \prod_{t \in \text{supp}(c)} \prod_{\substack{w \in \text{pos}(t): \\ r = d_G^t(w)}} \rho(r)^{c(t)} && \text{(commutativity)} \\
&= \left(\prod_{A \in N} \iota(A)^{c_G^I(A)} \right) \cdot \prod_{r \in R} \rho(r)^{\sum_{t \in \text{supp}(c)} \sum_{w \in \text{pos}(t): r = d_G^t(w)} c(t)} && (b^c \cdot b^d = b^{c+d}, 0^0 = 1, \text{ def. of } c_G^I) \\
&= \left(\prod_{A \in N} \iota(A)^{c_G^I(A)} \right) \cdot \prod_{r \in R} \rho(r)^{\sum_{t \in \text{supp}(c)} c(t) \cdot |\{w \in \text{pos}(t) \mid r = d_G^t(w)\}|} && \text{(distributivity)} \\
&= \left(\prod_{A \in N} \iota(A)^{c_G^I(A)} \right) \cdot \prod_{r \in R} \rho(r)^{c_G^R(r)} && \text{(def. of } c_G^R) \\
&= \frac{\prod_{A \in N} c_G^I(A)^{c_G^I(A)}}{\prod_{A \in N} |c|^{c_G^I(A)}} \cdot \frac{\prod_{r \in R} c_G^R(r)^{c_G^R(r)}}{\prod_{r \in R} c_G^N(\text{lhs}(r))^{c_G^R(r)}} && \text{(def. of } \iota \text{ and } \rho, \text{ comm., distr.)} \\
&= \frac{\prod_{A \in N} c_G^I(A)^{c_G^I(A)}}{|c|^{|c|}} \cdot \frac{\prod_{r \in R} c_G^R(r)^{c_G^R(r)}}{\prod_{A \in N} c_G^N(A)^{c_G^N(A)}} && (b^c \cdot b^d = b^{c+d}, \text{ def. of } c_G^I \text{ and } c_G^N, \text{ comm.})
\end{aligned}$$

This proves Equation 2.

B The function in Equation 5 is monotonically decreasing

We may transform Equation 5 as follows:

$$f(x, y) = \frac{x^x \cdot y^y}{(x+y)^{x+y}} = \left(\frac{x}{x+y} \right)^x \cdot \left(\frac{y}{x+y} \right)^y.$$

For positive arguments, the partial derivatives of f are

$$\begin{aligned}
\frac{\partial f(x, y)}{\partial x} &= \ln\left(\frac{x}{x+y}\right) \cdot \left(\frac{x}{x+y}\right)^x \cdot \left(\frac{y}{x+y}\right)^y, \text{ and} \\
\frac{\partial f(x, y)}{\partial y} &= \ln\left(\frac{y}{x+y}\right) \cdot \left(\frac{x}{x+y}\right)^x \cdot \left(\frac{y}{x+y}\right)^y.
\end{aligned}$$

For $x, y > 0$ the fractions are smaller than one. This means the logarithms are negative, hence the whole terms. So f is monotonically decreasing.

C Properties of mergers regarding bottom-up determinism

Lemma 1. *Let $G = (N, \Sigma, I, R)$ be an rtg, and let (\equiv_1) and (\equiv_2) be equivalence relations over N such that $\pi_{\equiv_1}(G)$ and $\pi_{\equiv_2}(G)$ are bottom-up deterministic. Let $(\equiv) = (\equiv_1) \cap (\equiv_2)$. Then also $\pi_{\equiv}(G)$ is bottom-up deterministic.*

Proof. Assume $\pi_{\equiv}(G)$ is not bottom-up deterministic. Then there are two rules $A_0 \rightarrow \sigma(A_1, \dots, A_{\text{rk}(\sigma)})$ and $B_0 \rightarrow \sigma(B_1, \dots, B_{\text{rk}(\sigma)})$ in R such that $A_i \equiv B_i$ for every $1 \leq i \leq \text{rk}(\sigma)$, but $A_0 \not\equiv B_0$. Hence, $A_i \equiv_1 B_i$ and $A_i \equiv_2 B_i$ for every $1 \leq i \leq \text{rk}(\sigma)$, and therefore $A_0 \equiv_1 B_0$ and $A_0 \equiv_2 B_0$. This implies $A_0 \equiv B_0$, which is a contradiction, so $\pi_{\equiv}(G)$ is bottom-up deterministic. q.e.d.

Lemma 2. *Let $G = (N, \Sigma, I, R)$ be an rtg, and let (\equiv) be an equivalence relation over N . Then there is a least (i.e. unique minimal) $(\hat{\equiv})$ such that $(\equiv) \subseteq (\hat{\equiv})$ and $\pi_{\hat{\equiv}}(G)$ is bottom-up deterministic.*

Proof. Existence: Consider (\equiv') such that $\forall A_1, A_2 \in N: A_1 \equiv' A_2$. Then (\equiv') satisfies the conditions.

Uniqueness: Assume there is a minimal $(\equiv') \neq (\hat{\equiv})$ satisfying the conditions. Then, by Lemma 1, $(\equiv') \cap (\hat{\equiv})$ would also satisfy the conditions, which contradicts that (\equiv') and $(\hat{\equiv})$ are minimal. Hence, $(\hat{\equiv})$ is unique. q.e.d.