Training PCFGs

Computational Linguistics

Alexander Koller

29 November 2019

Probabilistic CFGs

$S \rightarrow NP VP$	[1.0]	$VP \rightarrow V NP$	[0.5]
$NP \rightarrow Det N$	[0.8]	$VP \rightarrow VP PP$	[0.5]
$NP \rightarrow i$	[0.2]	$V \rightarrow shot$	[1.0]
$N \rightarrow N PP$	[0.4]	$PP \rightarrow P NP$	[1.0]
$N \rightarrow elephant$	[0.3]	$P \rightarrow in$	[1.0]
N → pyjamas	[0.3]	$\text{Det} \rightarrow \text{an}$	[0.5]
		$Det \rightarrow my$	[0.5]

(let's pretend for simplicity that Det = PRP\$)

Parse trees



"correct" = more probable parse tree

Today

- Parameters of PCFG = rule probabilities.
- How do we learn parameters from corpora?
 - maximum likelihood estimation
 - "hard EM" using Viterbi
 - "soft EM" using the inside-outside algorithm

ML Estimation

- Assume we have a treebank.
 - that is, every sentence annotated by hand with its "correct" parse tree
- Then we can use MLE to obtain rule probabilities:

$$P(A \to w) = \frac{C(A \to w)}{C(A \to \bullet)} = \frac{C(A \to w)}{\sum_{w'} C(A \to w')}$$

• Standard way of parameter estimation in practice. Works well, smoothing only needed for unknown words (or replace by POS tags).

Example



$N \rightarrow N PP$	[0]	$VP \rightarrow TV NP$	[1/4]
$N \rightarrow elephant$	[1/3]	$VP \rightarrow IV$	[1/4]
N → pyjamas	[2/3]	$VP \rightarrow VP PP$	[1/2]

Unsupervised estimation

- MLE works okay for English.
 - German: Tiger treebank exists, but is hard for PCFGs, e.g. because of free word order.
 - ▶ most other languages: phrase structure annotations unavailable, expensive to create → unsupervised methods?
- Unsupervised methods:
 - provide CFG, learn parameters from unannotated corpus
 - ▶ show first "hard EM", then "soft EM"
 - ideas instructive and generalize to other problems

"Hard" aka Viterbi EM

- In the absence of syntactic annotations, learner must invent its own parse trees.
- Viterbi EM:
 - start with some parameter estimate
 - produce "syntactic annotations" by computing best tree for each sentence using Viterbi
 - apply MLE to re-estimate parameters
 - repeat as long as needed
- This is *not* real EM!

Example



MLE on Viterbi parses



Some things to note

- In this example, the likelihood increased.
 - this need not always be the case for Viterbi EM
- Viterbi EM commits to a single parse tree per sentence. This has advantages and disadvantages:
 - parse tree easy to compute, and can simply apply MLE
 - ignores all uncertainty we had about correct parse (winning parse tree takes all)

Towards "real" (aka "soft") EM

idea: weighted counting of rules in all parse trees





Expected counts

• Define *expected count* of rule A → B C, based on previous parameter estimate.

$$E(A \to B \ C) = \sum_{t \in \mathcal{T}} P(t \mid w) \cdot C_t(A \to B \ C)$$

• If we have them, can re-estimate parameters:

$$P(A \to B \ C) = \frac{E(A \to B \ C)}{\sum_{r} E(A \to r)}$$

- Challenge: How to compute $E(A \rightarrow B C)$ efficiently?
 - we assume grammars in CNF here

Fundamental idea

Computing **µ**







Outside probabilities

 $\alpha(A, i, k) = \sum_{\substack{S \stackrel{d}{\Rightarrow}^{*} w_1 \dots w_{i-1} A w_k \dots w_n}} P(d)$

 $= \sum_{\substack{B \to A \ C \\ k < j \le n}} P(B \to A \ C) \cdot \beta(C, k, j) \cdot \alpha(B, i, j) + \sum_{\substack{B \to C \ A \\ 1 \le j < i}} P(B \to C \ A) \cdot \beta(C, j, i) \cdot \alpha(B, j, k)$



The Inside-Outside Algorithm

- Start with some initial estimate of parameters.
- For each sentence w, compute α , β , and μ .
- Compute expected counts $E(A \rightarrow B C)$.
 - sum expected counts over all sentences
 - remember that $P(w) = \beta(S, 1, n+1)$
- Re-estimate $P(A \rightarrow B C)$ from expected counts.
- Iterate until convergence.

Some remarks



- Inside-outside increases likelihood in each step.
- But huge problems with local maxima.
 - Carroll & Charniak 92 find 300 different local maxima for 300 different initial parameter estimates.
 - Improve by partially bracketing strings (Pereira & Schabes 92).
- Therefore, EM doesn't really work for totally unsupervised PCFG training.
- But extremely useful in refining existing grammars (Berkeley parser; see next time).

Summary

- Learning parameters of PCFGs:
 - maximum likelihood estimation from raw text
 - "hard EM": iterate MLE on Viterbi parses
 - EM: use inside-outside algorithm with expected rule counts
- PCFG parsing with MLE parse gets f-score in low 70's. Will improve on this next time (state of the art: 93).
- Have assumed that CFG is given and only parameters are to be learned. Will fix this later in this course.