

**10-710/11-763**  
**Fall 2013**  
**Homework 2**  
**Due: 10/6/2013**

Please submit your homework as a L<sup>A</sup>T<sub>E</sub>X-typeset PDF. Hand-drawn diagrams are ok. Direct clarification questions to [cmu-10710-2013-instructors@googlegroups.com](mailto:cmu-10710-2013-instructors@googlegroups.com)

## 1 Fun with HMMs

Given a finite label set  $\Omega$  and a finite observation alphabet  $\Sigma$ , a bigram HMM defines a joint distribution over label sequences  $\mathbf{y} = \langle y_1, y_2, \dots, y_n \rangle \in \Omega^n$  and observation sequences  $\mathbf{x} = \langle x_1, x_2, \dots, x_n \rangle \in \Sigma^n$  where  $n \geq 0$  as

$$p(\mathbf{x}, \mathbf{y}) = \left( \prod_{i=1}^n \gamma(y_i | y_{i-1}) \cdot \eta(x_i | y_i) \right) \cdot \gamma(\text{STOP} | y_n)$$

where  $y_0$  is assumed to be START.

1. Define an algorithm that runs in  $O(|\Omega|^2(n+k))$  time and generates  $k$  independent, random samples over the label variable  $y_i$ , drawn according to the distribution  $p(y_i | \mathbf{x})$ . This quantity can be derived from the HMM. (The inputs to the algorithm are the label set  $\Omega$ , the HMM parameters  $\gamma$  and  $\eta$ , the sequence  $\mathbf{x} = \langle x_1, x_2, \dots, x_n \rangle$ , and the position index  $i$ .) It's actually possible to do it with  $O(|\Omega|^2n + |\Omega|k)$  runtime, but that's not required.
2. Next, define an algorithm that runs in  $O(|\Omega|^2nk)$  time and generates  $k$  independent, random samples over the *label sequence* variable  $\mathbf{y}$ , drawn according to the conditional distribution  $p(\mathbf{y} | \mathbf{x})$ . (The inputs to the algorithm are the label set  $\Omega$ , the HMM parameters  $\gamma$  and  $\eta$ , and the sequence  $\mathbf{x}$ .)
3. Suppose I have an HMM such that  $\Omega = \{B, I, O\}$ . This model gives a joint distribution over word sequences bracketed with base noun phrases; this is also called "NP chunking." Given a structure  $(y_1^n, x_1^n)$  generated by this model, a "base noun phrase" corresponds to  $x_i^j = \langle x_i, x_{i+1}, \dots, x_j \rangle$  whenever  $(y_i^{j+1} = BI^{j-i}O) \vee (y_i^{j+1} = BI^{j-i}B) \vee ((j=n) \wedge (y_i^j = BI^{j-i}))$  (the first argument of the disjunction corresponds to a noun phrase followed by non-NP material, the second corresponds to a base NP that is followed by another base NP, and the third corresponds to a base NP at the end of the sentence). Define an algorithm that calculates the exact (marginal) probability that subsequence  $x_i^j$  is a base noun phrase, given the sequence and the HMM. (The inputs to the algorithm are the HMM parameters  $\gamma$  and  $\eta$ ,  $i$ ,  $j$ , and the sequence  $\mathbf{x}$ .)
4. Define an algorithm that runs in  $O(n)$  time and counts the number of valid NP chunkings (i.e., sequence labelings with  $\Omega$  as in item 1) for a sequence of length  $n$ . (The input to this algorithm is  $n$ .)

## 2 Fun with CRFs

HMMs define joint distributions over sequences of labels and tags. However, in many structured prediction applications, we are only interested in the conditional distribution of labels given observations,  $p(\mathbf{y} | \mathbf{x})$ . Conditional random fields (CRFs) permit this distribution to be modeled

directly, using the following parametric form:

$$p(\mathbf{y} \mid \mathbf{x}; \mathbf{w}) = \frac{\exp \mathbf{w}^\top \sum_{i=1}^n \mathbf{f}(x_i, y_i, y_{i-1})}{Z(\mathbf{x}; \mathbf{w})},$$

$$\text{where } Z(\mathbf{x}; \mathbf{w}) = \sum_{\mathbf{y}' \in \Omega^n} \exp \mathbf{w}^\top \sum_{i=1}^n \mathbf{f}(x_i, y'_i, y'_{i-1}),$$

$\mathbf{w} \in \mathbb{R}^d$  is parameter vector, and  $\mathbf{f} : \Sigma \times \Omega \times \Omega \rightarrow \mathbb{R}^d$  is a  $d$ -dimension feature vector function. Feature functions extract properties of local parts of the sequence that are useful for the prediction problem.

1. Define an algorithm that calculates the exact (marginal) probability that subsequence  $x_i^j$  is a base noun phrase, given the sequence and the CRF. (The inputs to the algorithm are the CRF parameters  $\mathbf{w}$  and the feature function  $\mathbf{f}$ ,  $i$ ,  $j$ , and the input sequence  $\mathbf{x}$ .)

To learn the parameters of a CRF, the conditional log likelihood of a sample of training data  $\mathcal{T} = \{\mathbf{x}_i^{(t)}, \mathbf{y}_i^{(t)}\}_{i=1}^n$ ,

$$\begin{aligned} \mathcal{L}(\mathbf{w}) &= \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{T}} \log p(\mathbf{y} \mid \mathbf{x}; \mathbf{w}) \\ &= \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{T}} \left[ \mathbf{w}^\top \sum_{i=1}^{|\mathbf{x}|} [\mathbf{f}(x_i, y_i, y_{i-1})] - \log Z(\mathbf{x}; \mathbf{w}) \right] \end{aligned} \quad (1)$$

is maximized by adjusting the parameters  $\mathbf{w}$ . There are a variety of ways to do this, but the most common ones involve computing the derivatives of the log likelihood function with respect to  $\mathbf{w}$ . Differentiating Eq. 1 using the chain rule for derivatives, we obtain

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{T}} \left[ \sum_{i=1}^{|\mathbf{x}|} \left[ \mathbf{f}(x_i, y_i, y_{i-1}) - \underbrace{\mathbb{E}_{p(\mathbf{y}' \mid \mathbf{x}; \mathbf{w})} \mathbf{f}(x_i, y'_i, y'_{i-1})}_{\text{Model expectation}} \right] \right].$$

2. Show how to obtain the “Model expectation” in the above equation by differentiating  $\log Z(\mathbf{x}; \mathbf{w})$  with respect to  $\mathbf{w}$ .
3. Give an algorithm that computes the expected value of  $\sum_{i=1}^{|\mathbf{x}|} [\mathbb{E}_{p(\mathbf{y}' \mid \mathbf{x}; \mathbf{w})} \mathbf{f}(x_i, y'_i, y'_{i-1})]$  for a training instance  $(\mathbf{x}, \mathbf{y})$ .

Let us now consider an alternative parameterization. Assume that each symbol  $\omega \in \Omega$  is associated with a vector  $\mathbf{v}_\omega \in \mathbb{R}^m$ , and that the feature vector function applies to a single symbol of  $\mathbf{f}(x_i) \in \mathbb{R}^d$ , and let:

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{\exp \sum_{i=1}^n \left[ \mathbf{f}(x_i)^\top \mathbf{W} \mathbf{v}_{y_i} + \mathbf{v}_{y_{i-1}}^\top \mathbf{V} \mathbf{v}_{y_i} \right]}{Z(\mathbf{x})},$$

where  $\mathbf{W} \in \mathbb{R}^{d \times m}$  and  $\mathbf{V} \in \mathbb{R}^{m \times m}$  are parameter matrices.

4. Given a set of training data give the derivatives of the log-likelihood with respect to  $\mathbf{W}$ ,  $\mathbf{V}$ , and  $\mathbf{v}_\omega$ .

### 3 Context Free Grammars

Following are a few rules<sup>1</sup> from a probabilistic context free grammar (PCFG) for English, along with their probabilities. Nonterminals and terminals are written in capital and small letters, respectively.

- $p(S \rightarrow NP VP) = 0.1, p(S \rightarrow VP) = 0.01$
- $p(NP \rightarrow N) = 0.1$
- $p(VP \rightarrow V) = 0.01, p(VP \rightarrow VP NP) = 0.1$
- $p(V \rightarrow \text{time}) = 0.0001, p(V \rightarrow \text{flies}) = 0.001$
- $p(N \rightarrow \text{time}) = 0.001, p(N \rightarrow \text{flies}) = 0.0001$

1. Find two complete parsing derivations of the sentence “time flies”. Which is more probable?
2. Why can’t the basic CYK<sup>2</sup> algorithm be used directly with this grammar?
3. Modify the rules of the grammar such that the CYK algorithm can be used.<sup>3</sup> Then, simulate, on paper, the CYK algorithm to find whether “time flies” can be parsed with this grammar. Show your work<sup>4</sup>.
4. The basic CYK algorithm solves the membership problem (i.e., whether a sentence belongs to the language defined with a CFG). There is an easy but inefficient way to modify CYK to also return the  $k$ -best parses.<sup>5</sup> Briefly and informally explain this easy modification.
5. The generalized CYK algorithm (Chappelier and Rajman, 1998) can be used to parse sentences with the given grammar without modifying the rules. Simulate the generalized CYK algorithm, on paper, with the given grammar and input sentence “time flies” showing your work. Does it give the same answer as part 3?
6. Features such as “nonterminal  $\mathbf{V}$  generates a terminal word that ends in  $\mathbf{ed}$ ” can help us encode our linguistic knowledge in a manageable way. The given PCFG uses categorical distributions. Briefly describe how a log-linear parameterization of those distributions allows us to use such features.

### 4 Parsing with Integer Linear Programming

A linear program in standard form looks like this:

$$\min_{\mathbf{x}} \mathbf{c}^T \mathbf{x} \text{ such that } \mathbf{A}\mathbf{x} \leq \mathbf{b}, \mathbf{x} \geq \mathbf{0}$$

Let  $N$  be the dimensionality of  $\mathbf{x}$  and let  $M$  be the number of linear constraints (not including the positivity constraints);  $\mathbf{A} \in \mathbb{R}^{M \times N}$  and  $\mathbf{b} \in \mathbb{R}^M$ . Any linear program can be transformed to this form.

<sup>1</sup>There are additional rules not shown; you do not need them for this problem.

<sup>2</sup>Some people call it CKY, others call it CYK. They are the same thing.

<sup>3</sup>You can use the algorithm described at <http://www.cs.nyu.edu/courses/fall107/V22.0453-001/cnf.pdf> – but it’s worth trying to work out how to do this manually first. The modified grammar should be equivalent to the original one (i.e. the languages defined by both grammars are identical).

<sup>4</sup>For example, you could show the CYK tables. This is not the only way to simulate CYK, though. So just pick a convenient way for you.

<sup>5</sup>It is probably easier to work it out on your own, but feel free to look at section 4.1 in (Huang and Chiang, 2005) The rest of section 4 in this paper describes less obvious and more efficient alternatives; we are not asking about these, although you should be aware they exist.

An *integer linear program* adds the constraint that each  $x_i$  is integer-valued. Often bounds constraints and integer constraints are merged into a constraint that each  $x_i \in \{0, 1\}$ .

1. Let  $\mathbf{w} = \langle w_0 = \$, w_1, w_2, \dots, w_n \rangle$  be a sentence. ( $\$$  is the root or “wall” symbol.) Let  $cost(w_i, w_j)$  be the cost associated with making  $w_i$  the parent of  $w_j$  in a dependency tree and let the cost of a dependency tree be

$$\sum_j cost(w_{\text{parent}(j)}, w_j)$$

Your job is to define an ILP such that minimizing  $\mathbf{x}$  can be converted back into a dependency tree such that:

- every word in  $\langle w_1, \dots, w_n \rangle$  has exactly one parent in  $\{w_0, w_1, \dots, w_n\}$ ;
- $w_0$  does not have a parent
- the solution to the ILP will correspond to the lowest-cost dependency tree.

To answer the question, you must define  $\mathbf{x}$  (in terms of the dependency tree),  $\mathbf{A}$ ,  $\mathbf{b}$ , and  $\mathbf{c}$ . Hint: you should be able to define a solution in which  $N = O(n^2)$  and  $M = O(n)$ . You don't need to convert to standard form, and it's okay to use equality and inequality constraints as long as they are linear.

2. Do part 1 again, but this time add and overall **projectivity** constraint on the tree:
  - if  $w_i$  is the parent of  $w_j$ , then  $\forall k \in (i, j) \cup (j, i)$ , the parent of  $w_k$  is also  $\in [i, j] \cup [j, i]$ . (Another equivalent statement is that for any  $i$ ,  $w_i$  and all its descendants for a contiguous substring.) Hint: this will require at least another  $O(n^2)$  linear constraints. Note: do not worry about cyclicity yet, unless you find it helpful to do so.
3. **Bonus question:** The dependency structures recovered above have not constrained to be *acyclic*. In practice we typically want dependency structures to be trees. Can you define constraints that will enforce acyclicity? You may introduce more variables, too, if you need to. You may also use non-linear constraints.

## References

- J.-C. Chappelier and M. Rajman. 1998. A generalized CYK algorithm for parsing stochastic CFG. In *Proc. of 1st Workshop on Tabulation in Parsing and Deduction*.
- L. Huang and D. Chiang. 2005. Better k-best parsing. In *In Proc. of the 9th International Workshop on Parsing Technologies*.