

Conditional Random Fields: An Introduction*

Hanna M. Wallach

February 24, 2004

1 Labeling Sequential Data

The task of assigning label sequences to a set of observation sequences arises in many fields, including bioinformatics, computational linguistics and speech recognition [6, 9, 12]. For example, consider the natural language processing task of labeling the words in a sentence with their corresponding part-of-speech (POS) tags. In this task, each word is labeled with a tag indicating its appropriate part of speech, resulting in annotated text, such as:

- (1) [PRP He] [VBZ reckons] [DT the] [JJ current] [NN account] [NN deficit] [MD will] [VB narrow] [TO to] [RB only] [# #] [CD 1.8] [CD billion] [IN in] [NNP September] [. .]

Labeling sentences in this way is a useful preprocessing step for higher natural language processing tasks: POS tags augment the information contained within words alone by explicitly indicating some of the structure inherent in language.

One of the most common methods for performing such labeling and segmentation tasks is that of employing hidden Markov models [13] (HMMs) or probabilistic finite-state automata to identify the most likely sequence of labels for the words in any given sentence. HMMs are a form of generative model, that defines a joint probability distribution $p(\mathbf{X}, \mathbf{Y})$ where \mathbf{X} and \mathbf{Y} are random variables respectively ranging over observation sequences and their corresponding label sequences. In order to define a joint distribution of this nature, generative models must enumerate all possible observation sequences – a task which, for most domains, is intractable unless observation elements are represented as isolated units, independent from the other elements in an observation sequence. More precisely, the observation element at any given instant in time may only directly

*University of Pennsylvania CIS Technical Report MS-CIS-04-21

depend on the state, or label, at that time. This is an appropriate assumption for a few simple data sets, however most real-world observation sequences are best represented in terms of multiple interacting features and long-range dependencies between observation elements.

This representation issue is one of the most fundamental problems when labeling sequential data. Clearly, a model that supports tractable inference is necessary, however a model that represents the data without making unwarranted independence assumptions is also desirable. One way of satisfying both these criteria is to use a model that defines a conditional probability $p(\mathbf{Y}|\mathbf{x})$ over label sequences given a particular observation sequence \mathbf{x} , rather than a joint distribution over both label and observation sequences. Conditional models are used to label a novel observation sequence \mathbf{x}_* by selecting the label sequence \mathbf{y}_* that maximizes the conditional probability $p(\mathbf{y}_*|\mathbf{x}_*)$. The conditional nature of such models means that no effort is wasted on modeling the observations, and one is free from having to make unwarranted independence assumptions about these sequences; arbitrary attributes of the observation data may be captured by the model, without the modeler having to worry about how these attributes are related.

Conditional random fields [8] (CRFs) are a probabilistic framework for labeling and segmenting sequential data, based on the conditional approach described in the previous paragraph. A CRF is a form of undirected graphical model that defines a single log-linear distribution over label sequences given a particular observation sequence. The primary advantage of CRFs over hidden Markov models is their conditional nature, resulting in the relaxation of the independence assumptions required by HMMs in order to ensure tractable inference. Additionally, CRFs avoid the label bias problem [8], a weakness exhibited by maximum entropy Markov models [9] (MEMMs) and other conditional Markov models based on directed graphical models. CRFs outperform both MEMMs and HMMs on a number of real-world sequence labeling tasks [8, 11, 15].

2 Undirected Graphical Models

A conditional random field may be viewed as an undirected graphical model, or Markov random field [3], globally conditioned on \mathbf{X} , the random variable representing observation sequences. Formally, we define $G = (V, E)$ to be an undirected graph such that there is a node $v \in V$ corresponding to each of the random variables representing an element Y_v of \mathbf{Y} . If each random variable Y_v obeys the Markov property with respect to G , then (\mathbf{Y}, \mathbf{X}) is a conditional random field. In theory the structure of graph G may be arbitrary, provided it represents the conditional independencies in the label sequences being modeled. However, when modeling sequences, the simplest and most common graph structure encountered is that in which the nodes corresponding to elements of

References

- [1] A. L. Berger. The improved iterative scaling algorithm: A gentle introduction, 1997.
- [2] A. L. Berger, S. A. Della Pietra, and V. J. Della Pietra. A maximum entropy approach to natural language processing. *Computational Linguistics*, 22(1):39–71, 1996.
- [3] P. Clifford. Markov random fields in statistics. In Geoffrey Grimmett and Dominic Welsh, editors, *Disorder in Physical Systems: A Volume in Honour of John M. Hammersley*, pages 19–32. Oxford University Press, 1990.
- [4] T. H. Cormen, C. E. Leiserson, and R. L. Rivest. *Introduction to Algorithms*. MIT Press/McGraw-Hill, 1990.
- [5] J. Darroch and D. Ratcliff. Generalized iterative scaling for log-linear models. *The Annals of Mathematical Statistics*, 43:1470–1480, 1972.
- [6] R. Durbin, S. Eddy, A. Krogh, and G. Mitchison. *Biological Sequence Analysis: Probabilistic Models of Proteins and Nucleic Acids*. Cambridge University Press, 1998.
- [7] E. T. Jaynes. Information theory and statistical mechanics. *The Physical Review*, 106(4):620–630, May 1957.
- [8] J. Lafferty, A. McCallum, and F. Pereira. Conditional random fields: probabilistic models for segmenting and labeling sequence data. In *International Conference on Machine Learning*, 2001.
- [9] A. McCallum, D. Freitag, and F. Pereira. Maximum entropy Markov models for information extraction and segmentation. In *International Conference on Machine Learning*, 2000.
- [10] S. Della Pietra, V. Della Pietra, and J. Lafferty. Inducing features of random fields. Technical Report CMU-CS-95-144, Carnegie Mellon University, 1995.
- [11] D. Pinto, A. McCallum, X. Wei, and W. B. Croft. Table extraction using conditional random fields. *Proceedings of the ACM SIGIR*, 2003.
- [12] L. Rabiner and B. H. Juang. *Fundamentals of Speech Recognition*. Prentice Hall Signal Processing Series. Prentice-Hall, Inc., 1993.
- [13] L. R. Rabiner. A tutorial on hidden Markov models and selected applications in speech recognition. *Proceedings of the IEEE*, 77(2):257–285, 1989.
- [14] A. Ratnaparkhi. A simple introduction to maximum entropy models for natural language processing. Technical Report 97-08, Institute for Research in Cognitive Science, University of Pennsylvania, 1997.

- [15] F. Sha and F. Pereira. Shallow parsing with conditional random fields. *Proceedings of Human Language Technology, NAACL 2003*, 2003.
- [16] C. E. Shannon. A mathematical theory of communication. *Bell System Tech. Journal*, 27:379–423 and 623–656, 1948.
- [17] H. M. Wallach. Efficient training of conditional random fields. Master’s thesis, University of Edinburgh, 2002.