

Undirected Graphical Models for Sequence Analysis

Fernando Pereira
University of Pennsylvania

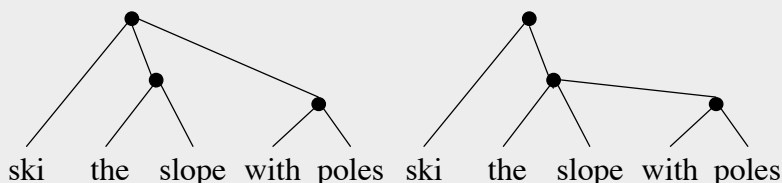
Joint work with John Lafferty,
Andrew McCallum, and Fei Sha

Analyzing Sequences

- Language
 - Syntactic structure
 - Sense tagging
 - Information extraction
- Biological sequences
 - Genes, regulatory regions
 - Secondary structure (folding)

Issues

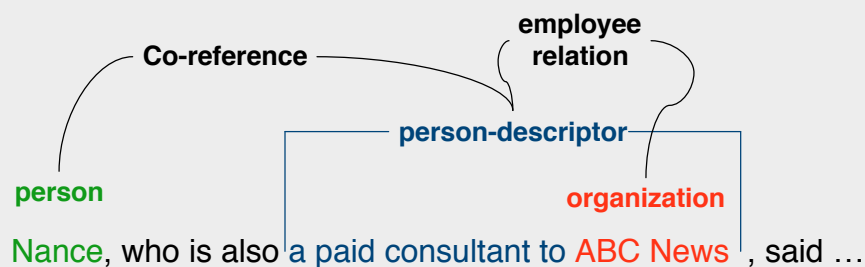
- Many interacting decisions



[ski [the slope] [with poles]] [ski [[the slope] [with poles]]]

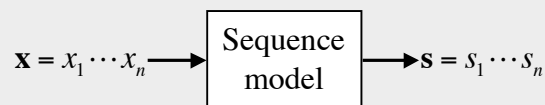
- A wide range of sequence features
- Computing an answer is relatively costly

Example: Information Extraction



- Information extraction as labeling:
 - States represent (parts of) entities
 - Relations
 - Cascaded labelers
 - Structured labels

Sequence Analysis as Labeling



- Labels (states) describe the role of corresponding inputs
- Appropriate for
 - Tagging
 - Shallow parsing
 - (Some) information extraction
 - Gene finding

Sequence Analysis

4

Local Labeling

- Train to minimize the per-decision loss in context

$$\hat{\square} = \arg \min_{\square} \prod_k \prod_{0 \leq i < |\mathbf{x}_k|} \text{loss}(s_{k,i} \mid \mathbf{x}_k, \mathbf{s}_k \setminus i; \square)$$

$$\mathbf{s}_k \setminus i = s_{k,1} \cdots s_{k,i-1} s_{k,i+1} \cdots s_{k,|\mathbf{x}_k|} \quad \boxed{\text{context of } s_{k,i}}$$

- Apply by *guessing* context and finding each lowest-loss label:

$$\hat{s}_i = \arg \min_{s_i} \text{loss}(s_i \mid \mathbf{o}, \hat{\mathbf{s}} \setminus i; \hat{\square})$$

Sequence Analysis

5

Global Labeling

- Minimize training labeling loss

$$\hat{\square} = \arg \min_{\square} \prod_k \text{Loss}(\mathbf{x}_k, \mathbf{s}_k \mid \square)$$

- Computing the best labeling:

$$\hat{\mathbf{s}} = \arg \min_{\mathbf{s}} \text{Loss}(\mathbf{x}, \mathbf{s} \mid \hat{\square})$$

- Efficient minimization requires:
 - A common currency for local labeling decisions
 - Efficient algorithm to combine the decisions

Method Tradeoffs

- *Assume*: Markovian label dependencies
- *Local*
 - ✓ Easy to use any classifier
 - ✗ No common currency for balancing labeling decisions
 - ✗ Search problem: guessing contexts
- *Global*
 - ✗ Restricted model forms
 - ✓ Easy to balance labeling decisions
 - ✓ Easy to compare alternative label sequences

A Global Model: HMMs

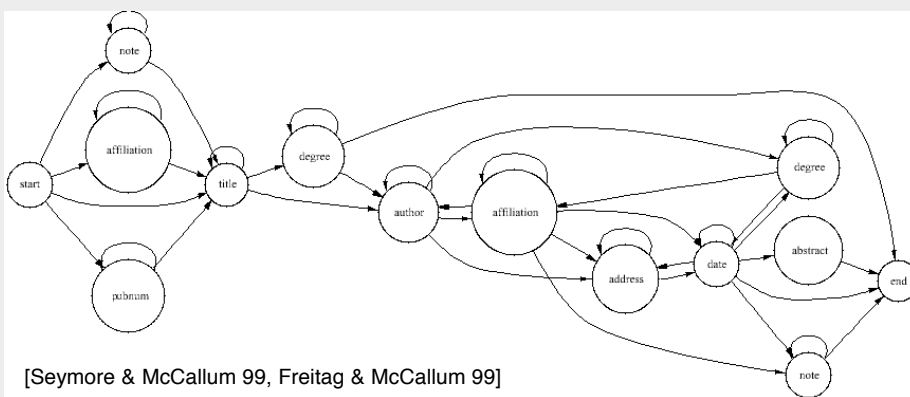
$$P(\mathbf{s}, \mathbf{x}) = P(s_0) \prod_i P(s_i | s_{i-1}) P(x_i | s_i)$$

- Train to maximize joint probability of input and label sequences
- Find most likely labeling for given input with Viterbi
- Commonly used: states emit sequences generated by a n -order Markov model

Sequence Analysis

8

Information Extraction with HMMs



- Input: text (word sequences)
- Label: which field does the word belong to

Sequence Analysis

9

Problems with HMMs

- Applications need richer input representation: multiple *overlapping* features, whole chunks of text

<i>Word features</i>	<i>Context features</i>
word identity	previous words
capitalization	next words
ends in "-tion"	markup
word in word list	starts sentence

- Generative models do not handle easily overlapping, non-independent features
- Alternative: *conditional* model $P(s|\mathbf{x})$

A Local Model: MEMMs

- Per-state conditional model

$$P(s \mid \mathbf{x} \mid i) = \frac{1}{Z(s, \mathbf{x} \mid i)} \exp \left[\mathbf{f}(s, s \mid \mathbf{x} \mid i) \right]$$

$\mathbf{x} \mid i \equiv$ input with indexing origin shifted to i

$$\text{Typical feature: } f(s, s \mid \mathbf{x}) = \begin{cases} 1 & \text{if } s = t \text{ \& } s \neq t \\ 0 & \text{otherwise} \end{cases} \quad \& \quad x_{\square 1} = x \ \& \quad x_0 = x[\square]$$

- Train each state model *separately*
- Decoding: Viterbi

Label Bias Problem

- Example (after Bottou '91):



- Bias toward states with fewer outgoing transitions.
- Per-state normalization does not allow the required $\text{score}(1, 2|ro) \ll \text{score}(1, 2|ri)$.
- *A challenge for all local models*

Sequence Analysis

12

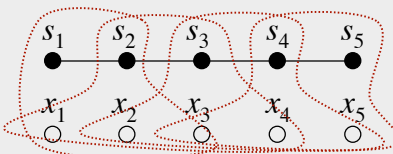
Conditional Undirected Models

- $P(\text{state sequence } s | \text{input sequence } x)$
instead of $P(x,s)$
- Allow arbitrary dependencies on x
- Efficient inference if dependencies within s are constrained
- States don't need to encode bounded dependency on past and future inputs

Sequence Analysis

13

Conditional Random Fields



- Markov on s , conditional dependency on \mathbf{x}

$$P_{\square}(\mathbf{s} | \mathbf{x}) = \frac{1}{Z_{\square}(\mathbf{x})} \prod_i \exp[\square \cdot \mathbf{f}(s_{i-1}, s_i, \mathbf{x}_{\square i})]$$

- Feature vector \mathbf{f} represents interactions between successive states and distribution of individual states given the input

From HMMs to CRFs

$$\mathbf{s} = s_1 \cdots s_n \quad \mathbf{x} = x_1 \cdots x_n$$

HMM
$$P(\mathbf{s} | \mathbf{x}) = \frac{P(s_0)}{P(\mathbf{x})} \prod_i P(s_i | s_{i-1}) P(x_i | s_i)$$

CRF
$$P(\mathbf{s} | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_i \exp \left[\begin{matrix} \square \square \cdot \mathbf{f}(s_i, s_{i-1}) \square \\ + \\ \square \square \cdot \mathbf{g}(s_i, x_i) \square \end{matrix} \right]$$

Decoding

- Use Viterbi algorithm to compute

$$\hat{\mathbf{s}} = \arg \max_{\mathbf{s}} P_{\square}(\mathbf{s} | \mathbf{x}) = \arg \max_{\mathbf{s}} \square \cdot \mathbf{F}(\mathbf{x}, \mathbf{s})$$

$$\mathbf{F}(\mathbf{x}, \mathbf{s}) = \prod_i \mathbf{f}(s_{i-1}, s_i, \mathbf{x} \square i)$$

- *Linear sequence models*: clean separation between
 - Decoding
 - Training
 - Generalizes linear classification

Efficient Estimation

- Matrix notation

$$M_i(s, s \square \mathbf{x}) = \exp \square \cdot \mathbf{f}(s, s \square \mathbf{x} \square i)$$

$$P_{\square}(\mathbf{s} | \mathbf{x}) = \frac{1}{Z_{\square}(\mathbf{o})} \prod_i M_i(s_{i-1}, s_i | \mathbf{x})$$

$$Z_{\square}(\mathbf{x}) = (M_1(\mathbf{x}) M_2(\mathbf{x}) \cdots M_{n+1}(\mathbf{x}))_{\text{start, stop}}$$

- Efficient normalization: *forward-backward* algorithm

Forward-Backward Calculations

- For any *path function* $G(\mathbf{s}) = \prod_i g_i(s_{i-1}, s_i)$

$$\begin{aligned}
 E_{\mathbf{S} \sim P(\mathbf{S}|\mathbf{x})} G(\mathbf{S}) &= \sum_{\mathbf{s}} P_{\square}(\mathbf{s} | \mathbf{x}) G(\mathbf{s}) \\
 &= \prod_i \frac{\alpha_i(\mathbf{x}) [g_{i+1} \beta_{i+1}(\mathbf{x})] \alpha_{i+1}^T(\mathbf{x})}{Z_{\square}(\mathbf{x})} \\
 \alpha_i(\mathbf{x}) &= \prod_{i-1} \alpha_{i-1}(\mathbf{x}) M_i(\mathbf{x}) \\
 \alpha_i^T(\mathbf{x}) &= M_{i+1}(\mathbf{x}) \alpha_i^T(\mathbf{x}) \\
 Z_{\square}(\mathbf{x}) &= \alpha_{n+1}(\text{end} | \mathbf{x}) = \alpha_0(\text{start} | \mathbf{x})
 \end{aligned}$$

- Easy generalization to constrained set of paths

Sequence Analysis

18

Training

- Maximize $L(\square) = \sum_k \log P_{\square}(\mathbf{s}_k | \mathbf{x}_k)$
- Log-likelihood *gradient*

$$\square L(\square) = \sum_k \left(\mathbf{F}(\mathbf{x}_k, \mathbf{s}_k) - E_{\mathbf{S} \sim P_{\square}(\mathbf{S}|\mathbf{x}_k)} \mathbf{F}(\mathbf{x}_k, \mathbf{S}) \right)$$

- Methods: iterative scaling, *conjugate gradient*, *L-BFGS*
- Partially-observable* case (labeled states)

$$\begin{aligned}
 P_{\square}(\mathbf{y} | \mathbf{x}) &= \sum_{\mathbf{s}: \ell(\mathbf{s})=\mathbf{y}} P_{\square}(\mathbf{s} | \mathbf{x}) \\
 \square L(\square) &= \sum_k \left[\sum_{\mathbf{s}: \ell(\mathbf{s})=\mathbf{y}_k} E_{\mathbf{S} \sim P_{\square}(\mathbf{S}|\mathbf{x}_k, \ell(\mathbf{S})=\mathbf{y}_k)} \mathbf{F}(\mathbf{x}_k, \mathbf{S}) - \sum_{\mathbf{s}} E_{\mathbf{S} \sim P_{\square}(\mathbf{S}|\mathbf{x}_k)} \mathbf{F}(\mathbf{x}_k, \mathbf{S}) \right]
 \end{aligned}$$

Sequence Analysis

19

Alternative Training Method

- Generalized perceptron [Collins 02]

$\mathbf{w}_0 = \mathbf{0}; j = 0$

for $t = 1, \dots, T$

 for $k = 1, \dots, N$

$\hat{\mathbf{s}} = \operatorname{argmax}_{\mathbf{s}} \mathbf{w}_j \cdot \mathbf{F}(\mathbf{x}_k, \mathbf{s})$

 if $\hat{\mathbf{s}} \neq \mathbf{s}_k$ then $\mathbf{w}_{j+1} = \mathbf{w}_j + \mathbf{F}(\mathbf{x}_k, \mathbf{s}_k) - \mathbf{F}(\mathbf{x}_k, \hat{\mathbf{s}})$ else $\mathbf{w}_{j+1} = \mathbf{w}_j$

$j \leftarrow j + 1$

$\bar{\mathbf{w}} = \frac{1}{NT} \sum_j \mathbf{w}_j$

- Questions:

- Could it lose something from Viterbi?
- Partially observable?

Sequence Analysis

20

Example: Shallow Parsing

- Noun phrase chunking from POS-tagged text

[Rockwell International Corp.] ['s Tulsa unit] said [it] signed [a tentative agreement] extending [its contract] with [Boeing Co.] to provide [structural parts] for [Boeing] ['s 747 jetliners]

- Standard benchmark task

Sequence Analysis

21

NP Chunking Results

Model	F
<i>24 SVM combination</i> [Kudo & Matsumoto 01]	94.22%
CRF [Sha & Pereira 03]	94.20%
Voted perceptron [Collins 02; Sha & Pereira 03]	94.09%
<i>Winnow</i> [Zhang, Damerau & Johnson 02]	93.89%
MEMM [Sha & Pereira 03]	93.70%

$$F = \frac{2PR}{P + R}$$

- **Warning:** different feature sets

Sequence Analysis

22

Fine Tuning

- Preconditioned conjugate gradient
 - Approximate diagonal of Hessian (exact diagonal is too expensive to compute)
- 820K features
 - Input predicate & transition predicate
 - Pre-compute input predicates
- Gaussian weight prior

Sequence Analysis

23

Further Questions

- Limited by dimensionality (number of features): kernels?
- Generalization bounds
- Parsing
 - Trees instead of chains
 - Inside-outside replaces forward-backward
 - Computationally challenging: large label set
- General graphs using loopy BP
 - Suggestive results for *collective classification* [Taskar & al 02]