
Efficient Inference on Sequence Segmentation Models

Sunita Sarawagi

IIT Bombay

SUNITA@IITB.AC.IN

Abstract

Sequence segmentation is a flexible and highly accurate mechanism for modeling several applications. Inference on segmentation models involves dynamic programming computations that in the worst case can be cubic in the length of a sequence. In contrast, typical sequence labeling models require linear time. We remove this limitation of segmentation models vis-a-vis sequential models by designing a succinct representation of potentials common across overlapping segments. We exploit such potentials to design efficient inference algorithms that are both analytically shown to have a lower complexity and empirically found to be comparable to sequential models for typical extraction tasks.

1. Introduction

Sequence segmentation models form the core of several applications including speech segmentation on phoneme boundaries (Keshet et al., 2005), information extraction (Sarawagi & Cohen, 2004), named entity recognition (McDonald et al., 2005), syntactic chunking (DauméIII & Marcu, 2005), shallow parsing, pitch accent prediction, and, protein/gene finding. Traditionally many of these applications have been artificially formulated as sequence labeling tasks at the expense of a loss of flexibility of features that can be exploited. This limitation is partly addressed by expanding the label set — for example, a popular choice in named entity recognition tasks (NER) is the Begin-Continue-End-Unique-other (BCEUO) encoding of entity labels (Borthwick et al., 1998), and in syntactic chunking tasks is the Begin-Inside-Outside (BIO) encoding of labels (Zhang et al., 2002). However, with the increasing diversity of settings where many of the applications are being deployed, it is imperative to exploit a richer variety of features than has been possible by sequence models. Examples of such segment level

features for extraction tasks are: the whole entity has an exact match in a database of entities, the length of the entity is between 4 and 8 words, the third or fourth token of the entity is a “-”, and the last three tokens in the entity are digits.

Sequence segmentation models provide a direct and natural way of encapsulating all such entity features. Given an input sequence $\mathbf{x} = x_1 \dots x_n$, a segmentation \mathbf{s} of \mathbf{x} consists of a sequence of variable length segments $\mathbf{s} = \langle s_1, \dots, s_p \rangle$ where each segment $s_j = \langle t_j, u_j, y_j \rangle$ consists of a *start position* t_j , an *end position* u_j , and a *label* $y_j \in Y$. Conceptually, a segment means that the tag y_j is given to all x_i 's between $i = t_j$ and $i = u_j$, inclusive. Each segment s_j can be associated with a vector of features that captures the dependence of its label on input properties in the neighborhood of the segment and the label of the segment before it. The goal during inference is to simultaneously find a segmentation of the input sequence and label each segment so as to maximize the total score over all segments.

The primary limitation of segmentation models compared to sequence models is the increased computational overheads of inference tasks (DauméIII & Marcu, 2005; Sarawagi & Cohen, 2004). The segmentation problem has been cast in a number of learning frameworks, including max-margin methods (McDonald et al., 2005) and CRF-based conditional likelihood methods (Sarawagi & Cohen, 2004). In all these different frameworks, there is a need to compute either the most likely segmentation as in cutting plane methods (Tsochantaridis et al., 2005) or various marginal probabilities as in likelihood methods and large margin methods optimized via exponentiated gradient algorithms (Bartlett et al., 2005). Both these inference problems are linear in the length n of the sequence for sequential labeling models, but can be quadratic (and is typically cube as we show in the paper) in n for segmentation models. A practical solution is to limit the maximum length of a segment by an a priori parameter L but even with such a choice, the training time for segmentation models has been observed to be 3–10 times worse than sequential models for extraction applications.

Appearing in *Proceedings of the 23rd International Conference on Machine Learning*, Pittsburgh, PA, 2006. Copyright 2006 by the author(s)/owner(s).

In this paper we seek to remove this limitation. The key insight we exploit is that most features are common across several overlapping segments. We propose a representation of features where such overlap can be compactly expressed. We then design inference algorithms where the running time is proportional to the number of features irrespective of how many segments a feature overlaps with. The algorithm pays the penalty of the segmentation models only when there are longer entity level features — when all the features are traditional token-level features the running time is the same as the standard Viterbi of sequence models. We achieve this without placing any hard limit on the length of a sequence. Empirical results show that with the new algorithm the running time for inference in segmentation models can be reduced by almost an order of magnitude. For typical information extraction models this makes the running time for training likelihood-based segmentation models comparable to that of sequential models while achieving higher accuracy because of the addition of a small, yet powerful set of entity-level features.

2. Graphical Models for Segmentation

We partition the potentials of a segmentation model into two types:

- Segment potentials $\theta_{i:j}(y)$ associated with a segment from i to j where all nodes from i to j have the base label y . Such potentials can be expressed in terms of features of an exponential model as $\theta_{i:j}(y) = \exp(\sum_{k=1}^K w_k f_k(\mathbf{x}, i, j, y))$ where w_k denotes the weight parameter of feature f_k .
- Transition potentials θ_i where an entry $\theta_i(y', y)$ denotes the potential for a segment starting at i getting a label y when the previous segment is labeled y' .

With these two types of potentials, the total score of a segmentation $\mathbf{s} = \langle s_1, \dots, s_p \rangle$ where each segment $s_j = \langle t_j, u_j, y_j \rangle$ consists of a start position t_j , an end position u_j , and a label $y_j \in Y$ can be expressed as $\prod_{j=1}^p \theta_{t_j}(y_{j-1}, y_j) \theta_{t_j:u_j}(y_j)$. During inference we need to find the sum or max marginal of a segment or edge potential. These inference problems can be solved optimally in polynomial time using an extension of the forward/backward algorithm for sequence labeling. The marginal probability of any segment potential can be computed via forward α and backward β messages which are expressed recursively as follows:

Let $\alpha_i(y)$ denote the sum of scores over all segmentations of the sequence between 1 to i where the last

segment has a label y .

$$\alpha_i(y) = \sum_{\max(i-L, 1) \leq i' \leq i} \sum_{y' \in \mathcal{Y}} \alpha_{i'-1}(y') \theta_{i'}(y', y) \theta_{i':i}(y) \quad (1)$$

where L , the maximum segment length is a parameter of the learning algorithm. The running time of this algorithm is $O(nL^2)$ where n is the input sequence length. This is because there are $O(nL)$ iterations and each segment $i' : i$ requires $O(i-i')$ work in computing the potential $\theta_{i':i}$.

Similarly, we can recursively compute the backward messages $\beta_i(y)$ that denotes the sum of scores over all segmentations of the sequence from $i+1$ to n with the segment ending at i having label y .

These can be used to find the marginal for a segment potential $i' : i$ labeled y as:

$$\frac{\sum_{y'} \alpha_{i'-1}(y') \theta_{i'}(y', y) \theta_{i':i}(y) \beta_i(y)}{Z(\mathbf{x})}$$

where $Z(\mathbf{x}) = \sum_y \alpha_n(y)$.

Thus, in $O(nL^2)$ time and two sets of messages of $O(n)$ length we can compute the marginal probability of any potential in the segmentation model. Similarly, we can find the most likely segmentation by replacing the two outer sums by two max terms in Equation 1.

Empirically, for typical NER tasks, segmentations models are found to require about 3–10 times the running time of sequential models. The value of the maximum length parameter L has to be chosen to be large enough to cover the largest segment because it is a hard limit on the length of the segment. Thus, for extraction from short text snippets, for example, citation records, L becomes comparable to n to cover long entities like “Titles”. This makes inference in segmentation models cubic in the length of the sequence. We next show how to address this shortcoming of segmentation models.

3. Efficient Inference Algorithms

A key insight we employ is that potentials for segmentations should be defined over a larger equivalence class of segments rather than for a single segment. We design a representation to express these equivalences succinctly through four kinds of potentials.

- $\psi_{i': \geq i}$ denotes potentials shared over all segments starting at i' but ending anywhere after or at i .
- $\psi_{\leq i' : i}$ denotes potentials shared over all segments ending at i but starting anywhere before or at i' .

- $\psi_{\leq i': \geq i}$ denotes potentials shared over all segments starting before or at i' and ending after or at i .
- $\psi_{i': i}$ denotes the usual full-segment potentials that apply to the exact segment between i' and i .

We show that in almost all applications of segmentations such kinds of potentials with varying levels of overlap across segments is commonplace.

Named Entity Recognition NER tasks are typically encoded via a BCEU encoding of labels. We can express all such potentials along with arbitrary other entity features without having to commit to a fixed encoding of labels. We list some examples of such potentials:

- First two tokens of segment starting at i are “Prof. Dr”: $\psi_{i: \geq i+2}$
- Last three tokens of segment ending at i are “of the ACM”: $\psi_{\leq i-3: i}$
- A segment of length l starting at i has high cosine similarity to a lexicon entry: $\psi_{i: i+l-1}$
- Segment contains “the” and “the” is the i -th word in \mathbf{x} : $\psi_{\leq i: \geq i}$
- The length of a segment starting at i is $> k$: $\psi_{i: \geq i+k}$

Syntactic chunking The set of features proposed in (DauméIII & Marcu, 2005) for syntactic chunking via segmentation also contains a mix of the four types of segment potentials. We highlight some examples:

- The token at position i before the start of a chunk is “it”: $\psi_{i+1: \geq i+1}$.
- The second and third token of a chunk starting at $i-1$ are stop-words: $\psi_{i-1: \geq i+1}$.
- The case pattern for a NP segment between tokens i and $i+2$ is three-caps: $\psi_{i: i+2}$.

Speech segmentation based on phoneme boundaries The input here is a sequence of frames and the goal is to segment the frames along a fixed number of phoneme boundaries. This can be treated as a segmentation problem and as shown in (Keshet et al., 2005) contains a variety of potential types as follows:

- Distance features between frames $i-s$ and $i+s$ for $s = 1, 2, 3, 4$ for a phoneme segment starting at i : $\psi_{i: \geq i+s}$
- The classification scores for the frames between i' and i being a phoneme: $\psi_{i': i}$

Our goal next is to exploit the succinct form of potentials to speed up the message passing algorithms

outlined in Section 2. Our new message passing algorithms can run in time that in the worst case is $O(nm + H)$ where m is the largest gap between the start and end boundary of potentials in *that sequence* and is typically smaller than the largest length of a segment and, H is the total number of potentials that are fired for a sequence. Thus, for sequence labeling tasks where $m = 2$, this will reduce to the standard $O(n)$ forward-backward algorithm, even though it can potentially output segments much larger than 2. We generalize this to the case of an arbitrary set of potentials. The main challenge in designing an efficient algorithm is that potentials could overlap in arbitrary ways and we cannot afford to pass messages only along transition edges as we did in Equation 1 with only full segment potentials.

We first show how this is done for forward and backward messages in Section 3.1. We then show in Section 3.2 how to directly compute marginals over potentials instead of summing over marginal probability of all segments the potential overlaps with.

3.1. Forward and Backward Terms

Let $\theta_{i': i}$ denote the product of all potentials applicable to segment $i' : i$; this can be expressed as $\theta_{i': i} = \prod \{\psi_{uv} : u = i' \vee u = (\leq k), k \geq i', v = i \vee v = (\geq j), j \leq i\}$. For example, in Figure 1 we show various potentials (as edges) for a sequence of length 9 where an arrow at any of the ends of an edge denotes potentials with open ends. Thus, the edge between nodes 3 and 5 denotes potential $\psi_{3: \geq 5}$ whereas the two edges between 4 and 7 denote the potential $\psi_{\leq 4: \geq 7}$ for the arrowed edge and $\psi_{4: 7}$ for the plain edge. For segment $4 : 7$ the total potential $\theta_{4: 7} = \psi_{4: \geq 5} \psi_{4: 7} \psi_{\leq 4: \geq 7} \psi_{\leq 6: 7}$

Equation 1 for computing α values with no L restriction can be written in matrix notation as:

$$\alpha_i = \sum_{i' \leq i} (\alpha_{i'-1} \theta_{i'}) * \theta_{i': i} \quad (2)$$

where the symbol “*” denotes element-wise multiplication of vectors of the same size. In the rest of the paper, we will drop the “*” to reduce clutter and assume it to be implicitly present when two vectors of the same length abut.

Our goal is to reuse computations performed for α_{i-1} to reduce the number of terms summed over for computing α_i . For this we design a method for decomposing the full segment potentials $\theta_{i': i}$ as $a'(i', i-1)a(i)$ where $a(i)$ is independent of i' and a' only involves potentials with the end boundary at $i-1$. We cannot hope to achieve this in general for all $i' \leq i-1$. Therefore, we define a function $A(i)$ such that for the

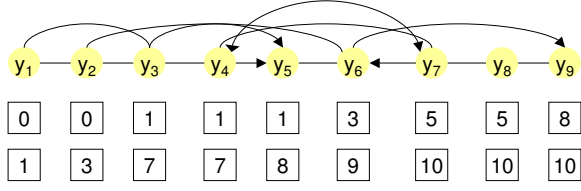


Figure 1. Potentials for a sequence of length 9. Potentials of the form $\psi_{j:\geq i}$ are represented with edges having an arrow at i , potentials of the form $\psi_{\leq j:i}$ are represented with edges having an arrow at j and so on. The first row of integers denotes $A(i)$ values for each i and the second row denotes the $B(i)$ values for each i . In this case $m=4$.

reduced set $i' \leq A(i-1)$ such a decomposition exists. We show how to design these functions.

Let $A(i)$ be the maximum index $j < i$ such that all potentials that start before or at j end before $i+1$.

$$A(i) = \max j : \forall \psi_{\ell:k} \quad \ell \leq j \Rightarrow k \leq i.$$

For the example in Figure 1 $A(6) = 3, A(7) = A(8) = 5$ and so on.

Let $\theta_{i':>i}$ denote the product of potentials common to all segments where the start boundary is i' and end boundary anywhere after i . Thus: $\theta_{i':>i-1} = \prod \{\psi_{r:s} : r = i' \vee r = (\leq j), j \geq i', s = (\geq k), k \leq i-1\}$.

Let $\theta_{i':(>i-1 \rightarrow i)}$ denote the product of potentials common to all segments that start at i' and end at i but not before i . That is, $\theta_{i':(>i-1 \rightarrow i)} = \prod \{\psi_{r:s} : r = i' \vee r = (\leq j), j \geq i', s = i \vee s = (\geq i)\}$. Note,

$$\theta_{i':i} = \theta_{i':>i-1} \theta_{i':(>i-1 \rightarrow i)} \quad (3)$$

From the definition of $A(i-1)$ we can claim that,

$$\theta_{i':(>i-1 \rightarrow i)} = \theta_{A(i-1):(>i-1 \rightarrow i)} \quad \text{if } i' \leq A(i-1) \quad (4)$$

Using equations 3 and 4, we can choose $a(i) = \theta_{A(i-1):(>i-1 \rightarrow i)}$ and $a'(i', i-1) = \theta_{i':>i-1}$ and use these to write the expression for α_i that reuses computations from α_{i-1} as follows: Let $\hat{\alpha}_i = \alpha_{i-1} \theta_i$ and $\hat{\alpha}_1 = 1$. Equation 2 can be rewritten as

$$\begin{aligned} \alpha_i &= \sum_{i' \leq A(i-1)} \hat{\alpha}_{i'} \theta_{i':>i-1} a_i + \sum_{A(i-1) < i' \leq i} \hat{\alpha}_{i'} \theta_{i':i} \\ &= \alpha_{i-1}^P a_i + \sum_{A(i-1) < i' \leq i} \hat{\alpha}_{i'} \theta_{i':i} \end{aligned} \quad (5)$$

where α_{i-1}^P denotes the sum of scores over all possible segmentations of sequence from 1 to i where the last segment's start boundary is before or at $A(i-1)$ and the end boundary after $i-1$. We compute α_i^P recursively as:

$$\alpha_i^P = \begin{cases} \alpha_{i-1}^P a_{>i} + \sum_{j=A(i-1)+1}^{A(i)} \hat{\alpha}_j \theta_{j:>i} & \text{if } A(i) > 0 \\ 0 & \text{if } A(i) \leq 0 \end{cases}$$

where $a_{>i}$ is like a_i except that we exclude potentials where the end boundary is strictly at i .

Thus, by maintaining an additional set of n forward terms denoting α_i^P we are able to compute α_i by summing over only $i - A(i-1)$ terms instead of $i-1$ terms. Note that $i - A(i-1) \leq m$, the maximum gap between the boundaries of any potential.

Example For the potentials in Figure 1 we show how to compute α_7 given $A(6) = 3$.

$$\begin{aligned} \alpha_7 &= \alpha_6^P a_7 + \hat{\alpha}_4 \theta_{4:7} + \hat{\alpha}_5 \theta_{5:7} + \hat{\alpha}_6 \theta_{6:7} \quad (\text{see Eq5}) \\ &= \alpha_6^P \psi_{\leq 4:\geq 7} \psi_{\leq 6:7} + \hat{\alpha}_4 \psi_{4:\geq 5} \psi_{4:7} \psi_{\leq 4:\geq 7} \psi_{\leq 6:7} + \\ &\quad \hat{\alpha}_5 \psi_{\leq 6:7} + \hat{\alpha}_6 \psi_{\leq 6:7} \end{aligned}$$

Similarly, with $A(7) = 5$ we compute α_7^P as

$$\alpha_7^P = \alpha_6^P \psi_{\leq 4:\geq 7} + \hat{\alpha}_4 \psi_{4:\geq 5} \psi_{\leq 4:\geq 7} + \hat{\alpha}_5$$

Beta terms The backward message β , for segmentation models is defined as

$$\beta_i = \theta_{i+1} \sum_{i' > i} \theta_{i+1:i'} \beta_{i'}$$

The computation of the beta messages can be optimized similarly to reuse terms from β_{i+1} in the computation of β_i . Accordingly, we define an index $B(i)$ such that for all $i' \geq B(i)$ we can decompose $\theta_{i:i'}$ as $b(i) b'(i+1, i')$ where $b(i)$ is independent of i' and b' involves potentials not before $i+1$.

Let $B(i)$ be the smallest index j such that all potentials that end after or at j do not have their start boundaries before i , that is,

$$B(i) = \min j : \forall \psi_{\ell:k} \quad k \geq j \Rightarrow \ell \geq i.$$

The last row in Figure 1 shows $B(\cdot)$ values. For example $B(4) = 7$ because there is no potential after position 7 that starts before position 4.

Let $\theta_{<i+1:i'}$ denote the product of all potentials in segments that end at i' but start anywhere to the left of $i+1$. That is, $\theta_{<i+1:i'} = \prod \{\psi_{r:s} : r = (\leq k), k \geq i+1, s = i' \vee s = (\geq j), j \leq i'\}$

Let $\theta_{<i+1 \rightarrow i:i'}$ denote the product of all potentials shared by segments ending at i' and starting at i but not after i . Thus, $\theta_{<i+1 \rightarrow i:i'} = \prod \{\psi_{r:s} : r = i \vee r = (\leq i), s = i' \vee s = (\geq j), j \leq i'\}$. Thus,

$$\theta_{i:i'} = \theta_{<i+1 \rightarrow i:i'} \theta_{<i+1:i'} \quad (6)$$

The definition of $B(i)$ implies that:

$$\theta_{<i+1 \rightarrow i:i'} = \theta_{<i+1 \rightarrow i:B(i+1)} \quad \text{if } i' \geq B(i+1) \quad (7)$$

Using equations 6 and 7 we can set $b(i) = \theta_{<i+1 \rightarrow i: B(i+1)}$ and $b'(i+1, i') = \theta_{<i+1: i'}$ to compute β_i without summing over all $n-i$ as follows:

$$\begin{aligned} \beta_i &= \theta_{i+1} \left(\sum_{B(i+2) > i' > i} \theta_{i+1: i'} \beta_{i'} + \sum_{i' \geq B(i+2)} b_{i+1} \theta_{<i+2: i'} \beta_{i'} \right) \\ &= \theta_{i+1} \left(\sum_{B(i+2) > i' > i} \theta_{i+1: i'} \beta_{i'} + b_{i+1} \beta_{i+2}^P \right) \end{aligned} \quad (8)$$

where $\beta_i^P = \sum_{i' \geq B(i)} \theta_{<i: i'} \beta_{i'}$, represents the sum over all segmentations \mathbf{s} where the first segment in \mathbf{s} starts before i and ends at or after position $B(i)$. β_i^P is computed recursively as:

$$\beta_i^P = \begin{cases} b_{<i} \beta_{i+1}^P + \sum_{j=B(i)}^{j=B(i+1)-1} \theta_{<i: j} \beta_j & \text{if } B(i) \leq n \\ 0 & \text{otherwise.} \end{cases}$$

where $b_{<i}$ is like b_i except that we exclude potentials where the start boundary is strictly at i .

Most likely segmentation The computation of the most likely segmentation can also be optimized via two dynamic programming equations similar to the two equations for α_i and α_i^P above so as to compute the maximum over at most m terms.

3.2. Computing Marginals around Potentials

The forward and backward messages can be combined to compute the marginals for various potentials. For normal segmentation models with potentials only of the form $\psi_{s:e}$, this involves only $O(1)$ computations of the form $\hat{\alpha}_s \theta_{s:e} \beta_e$. However, for potentials where the segment start and end boundaries are not fixed, computing the marginal for a potential of the form $\psi_{\leq s: \geq e}$ could require summing $O(n^2)$ such terms. We propose two ways to reduce the number of terms to be summed over. First, we use tricks like in the computation of α and β terms where we decompose potentials and depend on a parallel set of α^P and β^P terms. Second, we share computations across adjacent potentials. These two techniques together allow us to compute each marginal in $O(1)$ amortized time per potential. In the following sections we will use μ to denote un-normalized marginals i.e., the marginals before the division by $Z(\mathbf{x})$.

Potentials of the form $\psi_{\leq s: e}$ For such potentials the marginal $\mu_{\leq s: e}$ requires summing over s terms as follows:

$$\mu_{\leq s: e} = \sum_{i' \leq s} \hat{\alpha}_{i'} \theta_{i': e} \beta_e$$

We simplify this computation to reuse work across multiple related marginals as follows:

$$\mu_{\leq s: e} = \begin{cases} \alpha_{e-1}^P a_e \beta_e & \text{if } A(e-1) = s \\ \mu_{\leq (s-1): e} + \hat{\alpha}_s \theta_{s: e} \beta_e & \text{if } A(e-1) < s \end{cases}$$

The first case in the above equation follows from Equation 5 used to simplify the computation of α terms. The second case is a simple recursion and shows how we can reuse work in computing marginals of adjacent potentials. By the definition of $A(\cdot)$ in Eq: 4, we know that $s > A(e-1)$ and the case of $A(e-1) = s$ is a base case of the recursion.

Potentials of the form $\psi_{s: \geq e}$ In this case marginals require summing over $n-e$ terms as follows:

$$\mu_{s: \geq e} = \hat{\alpha}_s \sum_{i \geq e} \theta_{s: i} \beta_i$$

We simplify this computation so as to require summing over no more than m terms as follows:

$$\mu_{s: \geq e} = \begin{cases} \hat{\alpha}_s \beta_{s+1}^P b_s & \text{if } B(s+1) = e \\ \mu_{s: \geq (e+1)} + \hat{\alpha}_s \theta_{s: e} \beta_e & \text{if } B(s+1) > e \end{cases}$$

Edge potentials also fall in this class except that we need to restrict the previous label.

Potentials of the form $\psi_{\leq s: \geq e}$ In this case we need

$$\mu_{\leq s: \geq e} = \sum_{i' \leq s} \hat{\alpha}_{i'} \sum_{i \geq e} \theta_{i': i} \beta_i$$

We simplify this computation so as to not require summing over $O(n^2)$ terms as follows:

$$\mu_{\leq s: \geq e} = \begin{cases} \mu_{\leq s: \geq (e+1)} + \mu_{\leq s: e} & \text{if } B(s+1) > e \\ (\mu'(s-1) b_{<s} + \hat{\alpha}_s b_s) \beta_{s+1}^P & \text{if } B(s+1) = e \end{cases}$$

where $\mu'(s-1) = \mu_{\leq s-1: \geq B(s)} / \beta_s^P$. In the equation above the first case is obvious. The second case where $B(s+1) = e$ is derived next.

$$\begin{aligned} \mu_{\leq s: \geq e} &= \sum_{i' \leq s} \hat{\alpha}_{i'} \theta_{<s+1 \rightarrow i': e} \sum_{i \geq e} \theta_{<s+1: i} \beta_i \quad (\text{see Eq: 8}) \\ &= \sum_{i' \leq s} \hat{\alpha}_{i'} \theta_{<s+1 \rightarrow i': e} \beta_{s+1}^P \\ &= \sum_{i' \leq s} \hat{\alpha}_{i'} \theta_{<i'+1 \rightarrow i': e} \prod_{i' < j \leq s} \theta_{<j+1 \rightarrow <j: e} \beta_{s+1}^P \\ &= \sum_{i' \leq s} \hat{\alpha}_{i'} b_{i'} \prod_{i' < j \leq s} b_{<j} \beta_{s+1}^P \quad (\text{see Eq 7}) \\ &= ((\mu_{\leq s-1: \geq B(s)} / \beta_s^P) b_{<s} + \hat{\alpha}_s b_s) \beta_{s+1}^P \end{aligned}$$

It was tricky to get all the indices right in the above algorithm. We verified correctness by testing¹ that we get the same results on a direct computation of all the terms without such optimization.

¹Be careful about using the following code – I've only proven that it works, I haven't tested it. *Donald Knuth*

Algorithm 1 $\text{Alphas}(n, \psi_{i':\geq i}, \psi_{\leq i':i}, \psi_{\leq i':\geq i}, \psi_{i':i}, \theta_i)$

```

Initialize:  $\hat{\alpha}_1 = \mathbb{1}, A(0) = 0$ 
for  $i = 1 \dots n$  do
  Initialize  $\theta_{<i+1:i} = \mathbb{0}, \alpha_i = \mathbb{0}, \alpha_i^P = \mathbb{0}$ 
  for  $i' = i$  down to  $A(i-1) + 1$  do
    Get  $\theta_{<i':i}, \theta_{i':i}$  from  $\theta_{<i'+1:i}$  and  $\psi_s$ 
     $\alpha_i = \alpha_i + \hat{\alpha}_{i'} \theta_{i':i}$ 
  Compute  $A(i)$  from  $\psi_s$ 
  if  $A(i) > 0$  then
    Get  $\theta_{j:>i}, \theta_{<j:>i}$  from  $\theta_{<j:i}$  with  $j = A(i-1) + 1$ ,
    for  $j = A(i-1) + 1 \dots A(i)$  do
       $\alpha_i^P = \alpha_i^P + \hat{\alpha}_j \theta_{j:>i}$ 
      Get  $\theta_{j+1:>i}, \theta_{<j+1:>i}$  from  $\theta_{<j:>i}$  and  $\psi_s$ .
  if  $A(i-1) > 0$  then
    Get  $a_i, a_{>i}$  from  $a_{>i-1}$  and  $\psi_s$ .
     $\alpha_i = \alpha_i + \alpha_{i-1}^P a_i$ 
     $\alpha_i^P = \alpha_i^P + \alpha_{i-1}^P a_{>i}$ 
 $\hat{\alpha}_{i+1} = \alpha_i \theta_{i+1}$ 
    
```

3.3. Complexity Analysis

We show that the *worst case* complexity of inference is $O(nm + H)$ where m is the maximum span of any potential and H is the total number of features expressed as potentials. Consider the computation of α and α^P terms via Equation 5. The maximum number of terms summed over in each of the equations is m , the maximum gap between end positions of any potential. This explains the $O(nm)$ part. We explain the $O(H)$ part by showing that the θ -s can be computed in such a way that if the same potential is applicable to k adjacent segments the amount of work done is a constant independent of k . For this we start from $i' = i$ and decrease i' and in each iteration compute $\theta_{i':i}$ and $\theta_{<i':i}$ from a previous computation of $\theta_{<i'+1:i}$. We maintain the features in an efficient array-like structure such that for each $i' : i$ pair (there can be at most nm of these) and for each of the four possible kinds of potentials of the form $\psi_{\leq i':\geq i}, \psi_{\leq i':i}, \psi_{i':\geq i}, \psi_{i':i}$, we can retrieve all applicable features in $O(1)$ amortized time. Then, $\theta_{i':i}$ can be computed from $\theta_{<i'+1:i}$ by adding only the newly active features. We show how to compute the forward α terms via efficient potential reuse in Algorithm 1. The computation of betas is analogous.

In contrast, for the original algorithm the complexity in the *average case* is $O(nL + G)$ where L is a hard limit on the maximum segment length and is expected to be greater than m and G is the total number of features fired over all segments. We show that $O(G) = O(L^2 H)$ in the presence of token-level features. For example, consider a feature with the corresponding potential of the form $\psi_{\leq s:\geq s+\ell-1}$. This feature would be fired $(L - \ell + 1)(L - \ell + 2)/2$ times since it overlaps with that many segments. If all features were at the word-level with $\ell = 1$, then $O(G) = O(L^2 H)$. In tasks like

title/journal name extraction from citations where L is around 20 and several token-level features (like word and regular expression patterns indicators) are mixed with a few segment-level features (like match with a dictionary), this can lead to enormous savings as we see in the experimental section below. If we were to express complexity without involving the G and H terms and assumed that we need to perform $O(k)$ work to find potential of a k length segment, we get $O(nL^2)$ average case complexity for the old segmentation algorithm, which is reduced to $O(nm^2)$ worst case for the new algorithm.

The marginals μ can be computed in $O(mn)$ time as follows. All features are first sorted in increasing order of their start boundary followed by a decreasing order of their end boundary. The computation of μ s is done in the same order so that no storage needs to be allocated for them and the compute cost of μ is shared across all features with the same start and end boundary.

4. Experiments

We compare the optimized segmentation model (SegmentOpt) with the unoptimized model (Segment) and, for reference, also with the popular sequential labeling model with the begin-continue-end-unique (BCEU) encoding of labels (SequenceBCEU). We train each of these models to maximize the conditional likelihood objective regularized with a Gaussian prior on the parameters.

Datasets: We compared on three different extraction tasks spanning eight labels. The *Cora-Journals* dataset is the popular Cora citations dataset (Peng & McCallum, 2004) with two modifications: in the original Cora dataset, all authors are combined in a single label whereas we separate out individual authors, this makes the extraction problem harder but more meaningful and allows better exploitation of author name lexicons. We subset to include only journals (163 in all) to exploit journal and author names lexicons in the publicly available DBLP database. We extract title, authors and journal names. The *Address* corpus consists of 395 home addresses of students in a major university in a country with much less structure to their addresses than typical US addresses. We considered extraction of city names and state names from this corpus. For this dataset, we procured an external dictionary of city names and state names obtained from a postal database. The *Articles* dataset was created by starting from 400 journal entries over all the author's .bib files collected over several years. The unstructured records consisted of 100 citations obtained

Table 1. Comparing various methods on three IE tasks, with and without semi-Markov features. Column R is recall, P is precision and Time refers to the total training and testing time in seconds. SegmentOpt abbreviates SegmentOpt and Segment abbreviates Segment.

| dataset | Markov features ($m=1$) | | | | | Semi-Markov ($m = 7, 20, 20$) | | | | | |
|----------|---------------------------|------|------|------|------|---------------------------------|------|------|------|------|-----------|
| | method | R | P | F1 | Time | method | R | P | F1 | Time | $ G / H $ |
| Address | SequenceBCEU | 81.8 | 84.6 | 83.1 | 397 | Segment | 88.6 | 89.6 | 89.0 | 1130 | 8261 |
| | SegmentOpt | 85.5 | 87.9 | 86.5 | 292 | SegmentOpt | 88.4 | 89.9 | 89.1 | 342 | 2152 |
| Articles | SequenceBCEU | 81.7 | 84.8 | 82.9 | 220 | Segment | 85.5 | 86.4 | 85.8 | 2241 | 130826 |
| | SegmentOpt | 83.2 | 85.0 | 83.9 | 139 | SegmentOpt | 85.3 | 86.2 | 85.6 | 257 | 17015 |
| Cora | SequenceBCEU | 90.9 | 90.7 | 90.8 | 696 | Segment | 91.7 | 89.7 | 90.7 | 3623 | 193921 |
| | SegmentOpt | 91.4 | 89.9 | 90.6 | 557 | SegmentOpt | 91.6 | 89.7 | 90.6 | 800 | 21238 |

by searching Citeseer for citations of a random subset of authors appearing in the Bibtex database. The dataset is publicly available². We extracted individual author names, title and journal names. The main difference between this and the Cora dataset is that the lexicon has a greater overlap for this dataset.

Features: We used indicator features for the word itself and various surface patterns (capturing capitalization, digit pattern and delimiters) at the word and one position to its left and right. These features apply on the start, end and in-between part of segments – to get features equivalent to those in SequenceBCEU. Transition features are used with a history size of 1. These are all word-level features. The segment-level features we added were length feature for each possible segment-length value and for each available external lexicon we added features corresponding to the highest match with whole entities in the lexicon. We used two similarity functions (TFIDF, and JaroWinkler) which are known to work well for name-matching in data integration tasks (Cohen et al., 2003). For the benefit of the sequence models we also added word level dictionary match features using the same similarity functions. Thus, we believe that an equivalent set of features are available to the sequence and segmentation models.

We measure F1 accuracy of classifying the entire field correctly. The numbers are averaged over four random selections of 30% of the data for training and the rest for testing.

In Table 1 we compare SequenceBCEU with SegmentOpt (without any segment-level features), and Segment with SegmentOpt both with segment level features of maximum length 7 for the address data, and 20 for the other two. The numbers are averaged over all extracted labels in the dataset. We observe that F1 improves from 83 to 89 for the address data, 84 to 86 for Articles, and remains unchanged at 91 for the Cora dataset. This supports the conclusions in earlier work that segmentation models are better suited for

extraction tasks. The running time of Segment is a factor of 3–10 of SequenceBCEU. This shortcoming is removed with SegmentOpt where we are able to get the same accuracy advantage at the cost of only a factor 1–2 times higher running time compared to SequenceBCEU. Even with m of segment-level features as high as 20, the running time is comparable to the time of a sequence model. This gain is explained by the large difference in the sizes of features fired per sequence as shown in the last column. For example, the original set of 193,921 feature firings are represented equivalently as 21,238 feature firings (a factor of 10 reduction) with the succinct representation.

The improvement in running time achieved by SegmentOpt becomes more significant with increasing training fraction as we show in Figure 2 where running time for Segment increases sharply whereas SegmentOpt stays close to SequenceBCEU.

In Table 2 we illustrate the penalty of picking a wrong segment size L for Segment where we show F1 accuracy and running time with increasing segment size (parameter L for Segment and m for SegmentOpt). For a small segment size (10) in the Cora dataset accuracy of Segment is only 81 whereas with a segment size of 20 accuracy increases to 90.7, but this also doubles the running time. In contrast, for SegmentOpt accuracy remains steady around 90 and running time increases by only 15%.

Discussion In these experiments, the training objective was maximizing conditional likelihood; this requires marginal probabilities of potentials. We expect similar gains on max-margin methods trained via the exponentiated gradient optimization methods (Bartlett et al., 2005) since they require the same set of marginals. The training of max-margin methods via cutting plane algorithms (McDonald et al., 2005; Keshet et al., 2005; Tsochantaridis et al., 2005) that require finding the highest scoring segmentations will also show improvements, although we expect the gap to be smaller. Recently, DauméIII and Marcu (2005) propose a new structured learning framework where

²<http://www.it.iitb.ac.in/~sunita/data/personalBib.tar.gz>

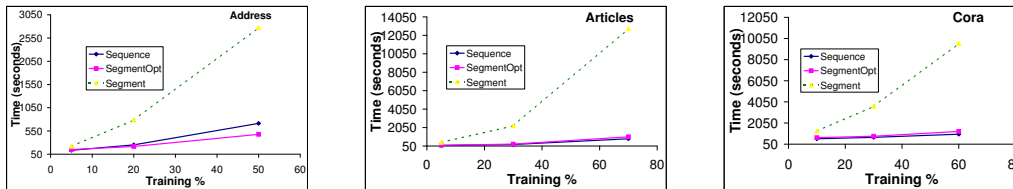


Figure 2. Running time against training set size for SequenceBCEU, Segment and SegmentOpt.

Table 2. Sensitiveness to segment size in the two semi-Markov models.

| | | Articles | | | | Cora | | | |
|---------|------------|----------------|------|------|------|------|------|------|------|
| | | Segment size → | | | | | | | |
| | | 10 | 15 | 20 | 25 | 10 | 15 | 20 | 25 |
| F1 | Segment | 83.5 | 85.0 | 85.8 | 85.8 | 81.0 | 88.8 | 90.7 | 90.3 |
| | SegmentOpt | 85.5 | 85.7 | 85.6 | 85.6 | 90.8 | 90.1 | 90.6 | 90.3 |
| Time(s) | Segment | 1026 | 1810 | 2241 | 2698 | 1687 | 2771 | 3623 | 3813 |
| | SegmentOpt | 175 | 245 | 257 | 260 | 712 | 767 | 800 | 781 |

the search for the best segmentation is integrated with parameter training. Our main idea of defining potentials over a larger equivalence class of segments and the modified dynamic programming algorithm is also applicable in this setting.

Concluding remarks Segmentation of sequences provide a natural, flexible, and high-accuracy modeling of several applications. The main limitation to their adoption vis-a-vis sequential models has been increased computational cost of inference algorithms. We removed this limitation by designing a compact feature representation for the mix of token-level and segment-level features in typical segmentation tasks, and exploited these to design efficient inference algorithms. The algorithm pays the penalty of entity-level features only when needed and is comparable to Viterbi on sequential models when entity features are a small fraction of the total feature set. In addition, there is no a priori hard limit on the maximum length of a segment. In future we would like to generalize to feature values that are continuous functions of the distance from the segment boundaries, for example, a feature that captures a Gaussian distribution on the length of the segment. Another interesting future work is generalizing to 2D segmentation models.

Acknowledgements The work reported here was supported by research grants from Microsoft Research and an IBM Faculty award.

References

- Bartlett, P. L., Collins, M., Taskar, B., & McAllester, D. (2005). Exponentiated gradient algorithms for large-margin structured classification. In L. K. Saul, Y. Weiss and L. Bottou (Eds.), *Advances in neural information processing systems 17*, 113–120. Cambridge, MA: MIT Press.
- Borthwick, A., Sterling, J., Agichtein, E., & Grishman, R. (1998). Exploiting diverse knowledge sources via maximum entropy in named entity recognition. *Sixth Workshop on Very Large Corpora New Brunswick, New Jersey. Association for Computational Linguistics.*
- Cohen, W. W., Ravikumar, P., & Fienberg, S. E. (2003). A comparison of string distance metrics for name-matching tasks. *Proceedings of the IJCAI-2003 Workshop on Information Integration on the Web (IIWeb-03)*. To appear.
- DauméIII, H., & Marcu, D. (2005). Learning as search optimization: approximate large margin methods for structured prediction. *ICML '05: Proceedings of the 22nd international conference on Machine learning* (pp. 169–176).
- Keshet, J., Shalev-Shwartz, S., & Singer, Y. (2005). Phoneme alignment using large margin techniques. *Workshop on the Advances in Structured Learning for Text and Speech Processing, NIPS*.
- McDonald, R., Crammer, K., & Pereira, F. (2005). Flexible text segmentation with structured multilabel classification. *Human Language Technology Conference Conference on Empirical Methods in Natural Language Processing (HLT/EMNLP)*.
- Peng, F., & McCallum, A. (2004). Accurate information extraction from research papers using conditional random fields. *HLT-NAACL* (pp. 329–336).
- Sarawagi, S., & Cohen, W. W. (2004). Semi-markov conditional random fields for information extraction. *NIPS*.
- Tsochantaridis, I., Joachims, T., Hofmann, T., & Altun, Y. (2005). Large margin methods for structured and interdependent output variables. *Journal of Machine Learning Research (JMLR)*, 6(Sep), 1453–1484.
- Zhang, T., Damerou, F., & Johnson, D. (2002). Text chunking based on a generalization of winnow. *J. Mach. Learn. Res.*, 2, 615–637.