

Towards Easier and Faster Sequence Labeling for Natural Language Processing: A Search-based Probabilistic Online Learning Framework (SAPO)

Xu Sun^{a,b,c,*}, Shuming Ma^{a,b}, Yi Zhang^{a,b}, Xuancheng Ren^{a,b}

^a*School of Electronics Engineering and Computer Science, Peking University*

^b*MOE Key Laboratory of Computational Linguistics, Peking University*

^c*Center for Data Science, Peking University*

Abstract

There are two major approaches for sequence labeling. One is the probabilistic gradient-based methods such as conditional random fields (CRF) and neural networks (e.g., RNN), which have high accuracy but drawbacks: slow training, and no support of search-based optimization (which is important in many cases). The other is the search-based learning methods such as structured perceptron and margin infused relaxed algorithm (MIRA), which have fast training but also drawbacks: low accuracy, no probabilistic information, and non-convergence in real-world tasks. We propose a novel and “easy” solution, a search-based probabilistic online learning method, to address most of those issues. The method is “easy”, because the optimization algorithm at the training stage is as simple as the decoding algorithm at the test stage. This method searches the output candidates, derives probabilities, and conducts efficient online learning. We show that this method, which is easy to implement, can support search-based optimization and obtain top accuracy with fast training and theoretical guarantee of convergence. Experiments on well-known tasks show that our method has better accuracy than CRF and BiLSTM¹.

*Corresponding author

Email addresses: xusun@pku.edu.cn (Xu Sun), shumingma@pku.edu.cn (Shuming Ma), zhangyi16@pku.edu.cn (Yi Zhang), renxc@pku.edu.cn (Xuancheng Ren)

¹The SAPO code is released at <https://github.com/lancopku/SAPO>

Keywords: natural language processing, sequence labeling, search-based learning, convergence

1. Introduction

Sequence labeling models are popularly used to solve structure dependent problems in a wide variety of application domains, including natural language processing, bioinformatics, speech recognition, and computer vision. To solve those problems, many sequence labeling methods have been developed, most of which are from two major categories. One is the probabilistic gradient-based learning methods such as conditional random fields (CRF) [16] and neural networks (e.g., RNN) [13]. The other is the search-based learning methods such as the margin infused relaxed algorithm (MIRA) [7] and structured perceptrons [5]. Other related work on sequence labeling includes maximum margin Markov networks [36] and structured support vector machines [37].

As for the probabilistic gradient-based learning methods such as CRF and RNN, they have high accuracy because of the exact computation of the gradient and probabilistic information. Nevertheless, those methods have critical drawbacks.

First, the probabilistic gradient-based methods typically do not support search-based optimization (search-based learning or decoding-based learning), which is important in sequence labeling problems with emphasis on the learning speed (e.g., for large-scale datasets). In the tasks with complex structures, the gradient computation is usually quite complicated and sometimes even intractable. This is mainly because dynamic programming for computing gradient is hard to scale for large-scale datasets. On the other hand, the search technique is easier to scale to large-scale datasets. This is because search-based learning is much simpler than gradient-based learning [29, 41, 30] — just search the promising output candidates and compare them with the oracle labels and update the weights accordingly.

Another category of sequence labeling methods is the search-based learning methods (decoding-based learning) such as structured perceptrons² and MIRA. A major advantage of those methods is that they support search-based learning such that the gradient is not needed and the learning is done by simply searching and comparing the promising output candidates with the oracle labels, and updating the model weights accordingly. As a by-product of the avoidance of gradient computation, those methods have faster training speed compared with probabilistic gradient-based learning methods such as CRF. However, there are also severe drawbacks of the existing search-based learning methods:

- First, the existing search-based learning methods such as structured perceptrons and MIRA have relatively low accuracy, compared with the probabilistic gradient-based learning methods such as CRF and RNN.
- Second, when applied to most of the real-world tasks, those search-based learning methods are non-convergent, i.e., they diverge in the training. As large margin classification models, theoretically those search-based learning methods have some convergent properties dependent on strict separability conditions. However, those strict separability conditions are not satisfiable in most real-world tasks, as demonstrated in many lines of prior work [31]. We will also show in the experiments that those search-based learning methods diverge dramatically as the training goes on, which makes the model accuracy go worse and worse.
- The existing search-based methods do not support probabilistic information. The magnitude of model weights grows dramatically as training

²The perceptron model used in the manuscript is the structured perceptron which includes features that incorporate transition features (structural information). The decoding of the structured perceptrons has to apply the Viterbi algorithm, because of the introduction of transition features. The decoding is the same with CRF models, which indeed considers the transition features. It is very different from the ordinary non-structured perceptron models, which can use greedy decoding because of the lack of transition features.

goes on, and there is no reliable probabilistic information that can be derived. We will also show the curves of the model weight magnitude in the experiments.

To address those issues, we propose a novel and “easy” solution, a search-based probabilistic online learning framework (SAPO), which can fix almost all of those drawbacks. The method is “easy” because the optimization algorithm at the training stage is as simple as the decoding algorithm at the test stage. The proposed method searches the top- n output candidates, derives probabilities based on the searched candidates, and conducts fast online learning by updating the model weights.

We show that the proposed method is of fast training speed which is comparable with structured perceptrons and MIRA, able to support search-based optimization and no need to calculate gradient, very easy to implement, with top accuracy which is even better than CRF and BiLSTM, and with theoretical guarantees of convergence towards the optimum given reasonable conditions. Experiments on well-known tasks show that our method has better accuracy than CRF and BiLSTM and almost as fast training speed as structured perceptrons and MIRA.

The contributions of this work are as follows:

- On the methodology side, we propose a general purpose search-based probabilistic online learning framework SAPO for sequence labeling. We show that SAPO can address a variety of issues of existing methods. Compared with probabilistic gradient-based learning methods such as CRF and RNN, the proposed method supports search-based learning such that avoiding complex gradient calculation, and with extra advantages on accuracy and training speed. Compared with search-based learning methods such as structured perceptron and MIRA, SAPO has much higher accuracy.
- For the proposed method, we provide theoretical and empirical justifications of convergence, even if the data is linearly non-separable. We provide

Algorithm 1 Search-based Probabilistic Online Learning Algorithm (SAPO)

- 1: **input:** top- n search parameter n , regularization strength λ , learning rate γ
 - 2: **repeat**
 - 3: Draw a sample $\mathbf{z} = (\mathbf{x}, \mathbf{y}^*)$ at random from training set S
 - 4: Based on \mathbf{w} , search the top- n outputs $Y_n = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n\}$
 - 5: For every $\mathbf{y}_k \in Y_n$, compute the probability $P_k = P(\mathbf{y}_k | \mathbf{x}, \mathbf{w})$
 - 6: For every $\mathbf{y}_k \in Y_n$, update the weights by $\mathbf{w} \leftarrow \mathbf{w} - \gamma P_k \mathbf{F}(\mathbf{x}, \mathbf{y}_k)$
 - 7: For \mathbf{y}^* , update the weights by $\mathbf{w} \leftarrow \mathbf{w} + \gamma \mathbf{F}(\mathbf{x}, \mathbf{y}^*)$
 - 8: Regularize the weights by $\mathbf{w} \leftarrow \mathbf{w} - \frac{\gamma \lambda}{|S|} \nabla R(\mathbf{w})$
 - 9: **until** Convergence
 - 10: **return** the learned weights \mathbf{w}^*
-

a novel theoretical analysis on the convergence of the proposed method, which does not rely on the assumption of linearly separable margin. On the other hand, structured perceptron and MIRA diverge in real-world tasks, which have linearly non-separable datasets, as to be shown in our experiments.

- On the application side, for several benchmark natural language processing tasks, including part-of-speech tagging, biomedical entity recognition, and phrase chunking, our simple search-based learning method can easily beat the strong baseline systems on those competitive tasks with faster speed.

2. Proposed method

We first describe the proposed search-based probabilistic online learning algorithm SAPO; then, we compare SAPO with existing methods.

2.1. Search-based probabilistic online learning

The proposed search-based probabilistic online learning algorithm SAPO has the key schemes as follows: top- n search (either exact search or approximate search), a scheme for calculating probabilities, perceptron-style update

for weights, and a regularizer on weights. We introduce the technical details of the key schemes as follows, after which we summarize the SAPO algorithm in Algorithm 1.

First, SAPO draws a training sample $\mathbf{z} = (\mathbf{x}, \mathbf{y}^*)$ at random from training set S , and searches for the top- n outputs:

$$Y_n = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n\}$$

In this work, each output \mathbf{y} is a structured label sequence $\mathbf{y} = \{y_1, y_2, \dots, y_l\}$, where l is the number of labels in the sequence. There are many methods to realize top- n search. One method uses the A* search algorithm [12]. An A* search algorithm with a Viterbi heuristic function can be used to produce top- n outputs one-by-one in an efficient manner. We use the backward Viterbi algorithm [42] to compute the admissible heuristic function for the forward-style A* search. This way, we can produce the top- n taggings efficiently.³

Then, for every $\mathbf{y}_k \in Y_n$, compute the probability in a log-linear fashion:

$$P_k \triangleq P(\mathbf{y}_k | \mathbf{x}, \mathbf{w}) \triangleq \frac{\exp[\mathbf{w}^T \mathbf{F}(\mathbf{x}, \mathbf{y}_k)]}{\sum_{\forall \mathbf{y} \in Y_n} \exp[\mathbf{w}^T \mathbf{F}(\mathbf{x}, \mathbf{y})]} \quad (1)$$

where \mathbf{w} is the vector of the model weights, $\mathbf{F}(\mathbf{x}, \mathbf{y}_k)$ is the feature vector based on \mathbf{x} and \mathbf{y}_k , and Y_n is simply the top- n outputs defined before. With this definition, we can see that $\sum_{k=1}^n P_k = 1$. That is, we use top- n search results to estimate the probability distribution, which is typically defined as:

$$P(\mathbf{y}_k | \mathbf{x}, \mathbf{w}) \triangleq \frac{\exp[\mathbf{w}^T \mathbf{F}(\mathbf{x}, \mathbf{y}_k)]}{\sum_{\forall \mathbf{y}} \exp[\mathbf{w}^T \mathbf{F}(\mathbf{x}, \mathbf{y})]} \quad (2)$$

As we can see, the only difference is the normalizer — we use top- n search results to estimate the normalizer. With the growth of n , this probability estimation in (1) goes more and more accurate towards the traditional probability in (2).

³Note that, although our search is “exact” top- n search, “exact” top- n search is not strictly required in the SAPO framework. In other words, we can replace exact A* search with non-exact beam search scheme for the SAPO algorithm. In the experiments we test both exact A* search and non-exact beam search with pruning (beam size is 50), and we find that there is almost no difference on the experimental results.

On the theoretical side, we will show in the theoretical analysis that this probability estimation can be arbitrarily-close to the traditional probability by using a proper n , and the SAPO algorithm is guaranteed to converge towards the optimum weights \mathbf{w}^* with an arbitrarily-close distance, given reasonable conditions. On the empirical side, we will show in experiments that the probability estimation is good enough for most real-world tasks even with $n = 5$ or $n = 10$.

After that, SAPO updates the weights in a perceptron fashion. For every $\mathbf{y}_k \in Y_n$, the weights are updated as follows:

$$\mathbf{w} \leftarrow \mathbf{w} - \gamma P_k \mathbf{F}(\mathbf{x}, \mathbf{y}_k) \quad (3)$$

As we can see, this is similar to the perceptron style update, except with an additional learning rate γ and a probabilistic scale P_k . On the other hand, for the oracle tagging \mathbf{y}^* , the weights are updated by

$$\mathbf{w} \leftarrow \mathbf{w} + \gamma \mathbf{F}(\mathbf{x}, \mathbf{y}^*) \quad (4)$$

As we can see, this is also similar to the perceptron style update. There is no need to use a probability scale here, because the probability is 1.

Finally, SAPO uses a weight regularizer with regularization strength λ , just like the stochastic regularization adopted in stochastic gradient descent (SGD) [2, 23, 31]. Following the regularization scheme of SGD, the regularization strength turns to $\lambda/|S|$ in the online learning setting [2, 23, 31]. Also, the regularization should be scaled with the learning rate γ . Thus, by using a regularizer denoted as $R(\mathbf{w})$, the regularization step is as follows:

$$\mathbf{w} \leftarrow \mathbf{w} - \frac{\gamma\lambda}{|S|} \nabla R(\mathbf{w}) \quad (5)$$

The regularizer $R(\mathbf{w})$ can be L_2 , L_1 , or other regularization terms. For simplicity, in this work we use the most widely used L_2 regularizer (a Gaussian prior).

The SAPO algorithm is summarized in Algorithm 1.

2.2. Comparison and discussion

Among the existing sequence labeling methods, the most similar and related methods to SAPO are structured perceptrons [5] and CRF [16].

If we compare SAPO with the structured perceptron [5] and CRF [16] with stochastic training, it is interesting to see that SAPO is like a “unification” of the structured perceptron and the stochastically trained CRF. The differences between CRF, structured perceptron, and SAPO mainly lie in how they estimate the parameters using the sequential information in the training phase. Structured perceptron updates the parameters only using the gold sequence path. CRF updates the parameters using an expectation of all the possible sequence paths. SAPO updates the parameters using the top- n possible sequence paths. If we neglect the learning rate and the regularizer term of SAPO, the structured perceptron algorithm [5] can be seen as an extreme case of SAPO with $n = 1$ (i.e., using top-1 search instead of top- n search). On the other hand, the stochastically trained CRF can be seen as another extreme case of SAPO with exponentially big n that enumerates over all possible output taggings (the only difference is that CRF uses dynamic programming instead of top- n search).

In other words, structured perceptron can be seen as SAPO with extremely small n , and CRF can be seen as SAPO with extremely big n . We argue that SAPO is more natural than both structured perceptrons and CRF — we should use a moderate value of n instead of an extremely small n (structured perceptrons) or an extremely huge n (CRF). As we will show in experiments and theoretical analysis, an extremely small n like structured perceptron will lead to low accuracy and non-convergent training, and an extremely large n like CRF will also lead to loss of accuracy (due to the overfitting of probabilities) and high computational cost. In practice, we find it good enough to use $n = 5$ or $n = 10$ for real-world tasks.

The MIRA algorithm also has a variation of Nbest MIRA which also uses top- n search [7]. Interestingly, it is also good enough to use $n = 5$ or $n = 10$ for Nbest MIRA [7, 19, 3]. Nevertheless, SAPO is substantially different compared with Nbest MIRA. The major difference is that SAPO has probability estimation

of different outputs while Nbest MIRA does not. Nbest MIRA treats different outputs equally without probability difference, which is why CRF cannot be seen as a special case of Nbest MIRA. Even if Nbest-MIRA uses extremely huge n in top- n search, it is not equivalent to CRF, and the difference is substantial. Also, there are other differences between SAPO and Nbest MIRA. For example, SAPO has the regularizer term and the learning rate and has no need to use the “minimum change” optimization criterion of MIRA during weight update.

3. Theoretical analysis

Here we give theoretical analysis on the objective function, update term, convergence conditions, and convergence rate.

3.1. Objective function and update term

Here we analyze the equivalent objective function of SAPO and the update term of SAPO. The SAPO algorithm (Algorithm 1) is a search-based optimization algorithm so that there is no need to compute the gradient of an objective function, and there is no explicit objective function used in the SAPO algorithm. Nevertheless, interestingly, we show that the SAPO algorithm is convergent and it converges towards the optimum weights \mathbf{w}^* which maximizes the objective function as follows:⁴

$$\text{maximize}_{\mathbf{w}} \sum_{i=1}^m \log P(\mathbf{y}_i^* | \mathbf{x}_i, \mathbf{w}) - \lambda R(\mathbf{w}) \quad (6)$$

where m is the number of training samples, i.e., $m = |S|$, and $R(\mathbf{w})$ is a weight regularization term for controlling overfitting. This objective function is similar to that of CRF. Equivalently, for the convenience of convex-based analysis, we denote the objective function $f(\mathbf{w})$ as the negative form of (6):

$$f(\mathbf{w}) = - \sum_{i=1}^m \log P(\mathbf{y}_i^* | \mathbf{x}_i, \mathbf{w}) + \lambda R(\mathbf{w}) \quad (7)$$

⁴The subscript of \mathbf{y} is overloaded here. For clarity throughout, \mathbf{y} with subscript i and usually with the $*$ mark refers to the tagging of the i 'th indexed training sample (e.g., \mathbf{y}_i^*), and \mathbf{y} with subscript k refers to the k 'th output of the search (e.g., \mathbf{y}_k).

We show that the SAPO algorithm converges towards the optimum \mathbf{w}^* which minimizes the convex objective function of $f(\mathbf{w})$:

$$\mathbf{w}^* = \underset{\mathbf{w}}{\text{minimize}} f(\mathbf{w}) \quad (8)$$

To clarify the theoretical analysis, we compare SAPO with the SGD (stochastic gradient descent) training scheme. Recall that the weight update has the following form in SGD [2, 23]:⁵

$$\mathbf{w} \leftarrow \mathbf{w} - \gamma \nabla f_{\mathbf{z}}(\mathbf{w}) \quad (9)$$

where $\nabla f_{\mathbf{z}}(\mathbf{w})$ is the stochastic gradient of $f(\mathbf{w})$ based on the sample \mathbf{z} , which has the following form when using the CRF objective function:

$$\begin{aligned} \nabla f_{\mathbf{z}}(\mathbf{w}) &= -\left\{ \mathbf{F}(\mathbf{x}, \mathbf{y}^*) - \sum_{\forall \mathbf{y}} P(\mathbf{y}|\mathbf{x}, \mathbf{w}) \mathbf{F}(\mathbf{x}, \mathbf{y}) - \frac{\lambda}{|S|} \nabla R(\mathbf{w}) \right\} \\ &= -\left\{ \mathbf{F}(\mathbf{x}, \mathbf{y}^*) - \sum_{\forall \mathbf{y}} \frac{\exp[\mathbf{w}^T \mathbf{F}(\mathbf{x}, \mathbf{y})]}{\sum_{\forall \mathbf{y}'} \exp[\mathbf{w}^T \mathbf{F}(\mathbf{x}, \mathbf{y}')] } \mathbf{F}(\mathbf{x}, \mathbf{y}) - \frac{\lambda}{|S|} \nabla R(\mathbf{w}) \right\} \end{aligned} \quad (10)$$

To make a comparison, we denote $\mathbf{s}_{\mathbf{z}}(\mathbf{w})$ as the (negative) SAPO update term for a sample \mathbf{z} such that

$$\mathbf{w} \leftarrow \mathbf{w} - \gamma \mathbf{s}_{\mathbf{z}}(\mathbf{w}) \quad (11)$$

Then, according to the procedure of SAPO algorithm, it is easy to check that $\mathbf{s}_{\mathbf{z}}(\mathbf{w})$ has the following form:

$$\begin{aligned} \mathbf{s}_{\mathbf{z}}(\mathbf{w}) &= -\left\{ \mathbf{F}(\mathbf{x}, \mathbf{y}^*) - \sum_{k=1}^n P_k \mathbf{F}(\mathbf{x}, \mathbf{y}_k) - \frac{\lambda}{|S|} \nabla R(\mathbf{w}) \right\} \\ &= -\left\{ \mathbf{F}(\mathbf{x}, \mathbf{y}^*) - \sum_{\forall \mathbf{y} \in Y_n} \frac{\exp[\mathbf{w}^T \mathbf{F}(\mathbf{x}, \mathbf{y})]}{\sum_{\forall \mathbf{y}' \in Y_n} \exp[\mathbf{w}^T \mathbf{F}(\mathbf{x}, \mathbf{y}')] } \mathbf{F}(\mathbf{x}, \mathbf{y}) - \frac{\lambda}{|S|} \nabla R(\mathbf{w}) \right\} \end{aligned} \quad (12)$$

As we can see from (10) and (12), by increasing n , the SAPO update term $\mathbf{s}_{\mathbf{z}}(\mathbf{w})$ can be arbitrarily-close to the stochastic gradient $\nabla f_{\mathbf{z}}(\mathbf{w})$. More formally, we

⁵In practice, SGD and SAPO can use decayed learning rate or fixed learning rate. Following [23, 30], for the convenience of theoretical analysis, our theoretical analysis is more focused on SGD and SAPO with fixed learning rate.

define

$$\delta_{\mathbf{z}}(\mathbf{w}) = \nabla f_{\mathbf{z}}(\mathbf{w}) - \mathbf{s}_{\mathbf{z}}(\mathbf{w}) \quad (13)$$

Then, for any $\epsilon \geq 0$, there is at least a corresponding n such that,

$$\delta_{\mathbf{z}}(\mathbf{w}) \leq \epsilon \quad (14)$$

In other words, when n is increasing, the approximation is expected to be more and more accurate and finally reaches the point where $\delta_{\mathbf{z}}(\mathbf{w}) \leq \epsilon$.

3.2. Optimum, convergence, and convergence rate

Recall that $f(\mathbf{w})$ is the sequence labeling objective function, and $\mathbf{w} \in \mathcal{W}$ is the weight vector. Taking the time stamp t into consideration, the SAPO update (11) can be reformulated as follows:

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \gamma \mathbf{s}_{\mathbf{z}_t}(\mathbf{w}_t) \quad (15)$$

To state our convergence analysis results, we make several assumptions following [22]. We assume f is strongly convex with modulus c , that is, $\forall \mathbf{w}, \mathbf{w}' \in \mathcal{W}$,

$$f(\mathbf{w}') \geq f(\mathbf{w}) + (\mathbf{w}' - \mathbf{w})^T \nabla f(\mathbf{w}) + \frac{c}{2} \|\mathbf{w}' - \mathbf{w}\|^2 \quad (16)$$

where $\|\cdot\|$ means 2-norm $\|\cdot\|_2$ by default in this work. When f is strongly convex, there is a global optimum/minimizer \mathbf{w}^* . We also assume Lipschitz continuous differentiability of ∇f with the constant q , that is, $\forall \mathbf{w}, \mathbf{w}' \in \mathcal{W}$,

$$\|\nabla f(\mathbf{w}') - \nabla f(\mathbf{w})\| \leq q \|\mathbf{w}' - \mathbf{w}\| \quad (17)$$

Also, let the norm of $\mathbf{s}_{\mathbf{z}}(\mathbf{w})$ be bounded by $\kappa \in \mathbb{R}^+$:

$$\|\mathbf{s}_{\mathbf{z}}(\mathbf{w})\| \leq \kappa \quad (18)$$

Moreover, it is reasonable to assume

$$\gamma c < 1 \quad (19)$$

because even the ordinary gradient descent methods will diverge if $\gamma c > 1$ [23].

Based on the assumptions, we show that SAPO converges towards the minimum \mathbf{w}^* of $f(\mathbf{w})$ with an arbitrary-close distance, and the convergence rate is given as follows.

Theorem 1 (Optimum, convergence, and rate). *With the conditions (16), (17), (18), (19), let $\epsilon > 0$ be a target degree of convergence. Let τ be an approximation-based bound from $\mathbf{s}(\mathbf{w})$ to $\nabla f(\mathbf{w})$ such that*

$$[\nabla f(\mathbf{w}) - \mathbf{s}(\mathbf{w})]^T (\mathbf{w} - \mathbf{w}^*) \leq \tau \quad (20)$$

where \mathbf{w} is a historical weight vector that updated during SAPO training, and $\mathbf{s}(\mathbf{w})$ is expected $\mathbf{s}_z(\mathbf{w})$ over \mathbf{z} such that $\mathbf{s}(\mathbf{w}) = \mathbb{E}_z[\mathbf{s}_z(\mathbf{w})]$. Since $\mathbf{s}(\mathbf{w})$ can be arbitrary-close to $\nabla f(\mathbf{w})$ by increasing n , SAPO can use the smallest n as far as the following holds:

$$\tau \leq \frac{c\epsilon}{2q} \quad (21)$$

Let γ be a learning rate as

$$\gamma = \frac{c\epsilon - 2\tau q}{\beta q \kappa^2} \quad (22)$$

where we can set β to be any value as far as $\beta \geq 1$. Let t be the smallest integer satisfying

$$t \geq \frac{\beta q \kappa^2 \log(qa_0/\epsilon)}{c(c\epsilon - 2\tau q)} \quad (23)$$

where a_0 is the initial distance such that $a_0 = \|\mathbf{w}_0 - \mathbf{w}^*\|^2$. Then, after t updates of \mathbf{w} , SAPO converges towards the optimum such that

$$\mathbb{E}[f(\mathbf{w}_t) - f(\mathbf{w}^*)] \leq \epsilon \quad (24)$$

The proof is in Appendix A.

This theorem shows that the approximation based learning like SAPO is also convergent towards the optimum of the objective function. Thus, we can approximate the true gradient by top- n search and still keep the convergence properties, without having to calculate exact gradients such as the training of CRF.

More specifically, the theorem shows that SAPO is able to converge towards the optimum of the objective function with arbitrarily close distance ϵ , as far as the SAPO update term $\mathbf{s}(\mathbf{w})$ is a ‘‘close-enough approximation’’ (i.e., satisfying (21)) of the true gradient $\nabla f(\mathbf{w})$. Since $\mathbf{s}(\mathbf{w})$ can be arbitrary-close to $\nabla f(\mathbf{w})$ by increasing n , SAPO can use the smallest n as far as the close-enough

approximation (21) is satisfied. In practice, we find that setting n to 5 or 10 already empirically satisfies the close-enough approximation in most of the real-world tasks. Moreover, the convergence rate is given in the theorem — SAPO is guaranteed to converge with t updates, and t is the smallest integer satisfying (23).

This analysis also explains why the structured perceptron algorithm [10, 5] does not converge in most of the practical tasks. As discussed before, the structured perceptron algorithm can be essentially treated as an extreme case of SAPO, which uses an extremely small n as 1. In most cases, the use of $n = 1$ does not satisfy the close-enough approximation condition of (21). Thus in most cases the structured perceptron algorithm has a bad approximation over the true gradient and it diverges (as we will show in experiments).

As for the real-world tasks, the datasets are often linearly non-separable, so the structured perceptron and MIRA will diverge. However, according to the theoretical analysis, SAPO is able to remain convergent even when the data is linearly non-separable. Our experiments in Section 4 will show how SAPO is still convergent in the real-world tasks.

4. Experiments

We describe the real-world tasks for the experiments, the experimental settings, and the experimental results as follows.

4.1. Tasks

We conduct experiments on natural language processing tasks with quite diverse characteristics. The natural language processing tasks include (1) part-of-speech tagging, (2) biomedical named entity recognition, and (3) phrase chunking. All of the tasks use boolean features. From tasks (1) to (3), the average length of samples (i.e., the number of tags per sample) is quite different, being 23.9, 26.5, 46.6, respectively. The dimension of tags $|\mathcal{Y}|$ is also very diversified among tasks, with $|\mathcal{Y}|$ ranging from 5 to 45.

Part-of-Speech Tagging (POS-Tag): Part-of-Speech (POS) tagging is an important and highly competitive task in natural language processing. We use the standard benchmark dataset in prior work [5, 40], which is derived from PennTreeBank corpus and uses sections 0 to 18 of the Wall Street Journal (WSJ) for training (38,219 samples), and sections 22-24 for testing (5,462 samples). Following the previous work [39], we use features based on unigrams and bigrams of neighboring words, and lexical patterns of the current word, with 393,741 raw features⁶ in total. Following previous work, the evaluation metric for this task is per-word accuracy.

Biomedical Named Entity Recognition (Bio-NER): This task is from the *BioNLP-2004 shared task*, which is to recognize 5 kinds of biomedical named entities (*DNA*, *RNA*, etc.) on the *MEDLINE* biomedical text corpus [20, 4]. There are 17,484 training samples and 3,856 test samples. Following the previous work [39], we use word pattern features and POS features, with 403,192 raw features in total. The evaluation metric is balanced F-score.

Phrase Chunking (Chunking): In the phrase chunking task, the non-recursive cores of noun phrases called base NPs are identified. The phrase chunking data is extracted from the data of the *CoNLL-2000 shallow-parsing shared task* [27]. The training set consists of 8,936 sentences, and the test set consists of 2,012 sentences. Following the previous work [39], we use the feature templates based on word n-grams and part-of-speech n-grams, with 264,818 raw features in total. Following previous studies, the evaluation metric for this task is balanced F-score.

4.2. Experimental settings

We compared the proposed SAPO algorithm with strong baselines in the existing literature, including both probabilistic gradient-based learning methods and search-based learning methods. For the probabilistic gradient-based learn-

⁶Raw features are those observation features based only on \mathbf{x} , i.e., no combination with tag information.

ing methods, we choose the arguably most popular model CRF [16], bidirectional LSTM (Bi-LSTM) [28], and bidirectional LSTM CRF (Bi-LSTM-CRF) [14] as the baselines. The CRF is with the widely used L_2 regularization and is trained with the standard SGD training algorithm. The Bi-LSTM and Bi-LSTM-CRF are trained with Adam optimizing algorithm [15].

For search-based learning methods, we choose structured perceptrons (Perc) [5] and MIRA [7], which are arguably the most popular search-based learning methods, as the baselines. In most cases, the averaged versions of structured perceptrons and MIRA work empirically better than naive versions of structured perceptron and MIRA [5, 19, 8, 3]. Thus we also compare SAPO with averaged versions of structured perceptrons and MIRA. To differentiate the naive and averaged versions, we denote them as Perc-Naive, Perc-Avg, MIRA-Naive, MIRA-Avg, respectively. Moreover, the MIRA method has the Nbest versions [7, 19], which adopt top- n search and update instead of Viterbi search and update. We also choose Nbest versions of MIRA as the additional baselines. We denote the Nbest MIRA with naive training as MIRA-Nbest-Naive and denote the one with averaged training as MIRA-Nbest-Avg.

The regularization strength λ of CRF is tuned among values 0.1, 0.5, 1, 2, 5 and are determined on the development data provided by the standard dataset (POS-Tag) or simply via 4-fold cross validation on the training set (Bio-NER and Chunking). With this automatic tuning for regularization strength, we set it to be 2, 5, 1 for POS-Tag, Bio-NER, and Chunking tasks, respectively. To give no tuning advantage to SAPO, SAPO simply uses the same regularizer and the same learning rate as CRF does. All the tuning is based on CRF, and there is no additional tuning for SAPO.

Also, the proposed SAPO algorithm uses the same top- n search scheme as the Nbest MIRA does. As shown in the previous work [7, 19, 3], it is good enough to use $n = 5$ or $n = 10$ for Nbest MIRA. It is also good enough to use $n = 5$ or $n = 10$ for the proposed SAPO algorithm. Thus, we set $n = 5$ for Nbest MIRA and SAPO for fast speed.

The features used are based on previous work [39]. The features used for

chunking are unigrams and bigrams of neighboring words, as well as unigrams, bigrams and trigrams of neighboring POS tags. The features used for Bio-NER are unigrams of neighboring chunk tags, substrings (shorter than 10 characters) of the current word, and the morphological features of the word, as well as the features used in the chunking experiments. For the features of POS-Tag, we use unigrams and bigrams of neighboring words, prefixes and suffixes of the current word, and some characteristics of the word. We also normalize the current word by turning all capital letters into lower case and converting all the numerals into ‘#’, and used the normalized word as a feature. Those features are exactly based on the previous work [39]. For the neural models, the uni-gram features in the feature set use pre-trained Senna word embeddings [6] with 50 dimensions for initialization. The embeddings are optimized by backpropagation during training.

The labeling scheme is the same for all models (including neural models). For POS-Tag, it is the original POS tag. For BioNER and Chunking, it is the BIO scheme following the previous work [39].

All Experiments are performed on a computer with the Intel(R) Xeon(R) 3.0GHz CPU. For fair comparison, we set the batch size to 1 for all experiments.

4.3. Experimental results

The experimental results in terms of accuracy/F-score and the computation cost are shown in Figure 1, Figure 2, Figure 3, and Figure 4. As we can see, although the tasks involve diverse feature types and different characteristics, the results are quite consistent — the proposed SAPO algorithm has the best accuracies/F-scores in all of the three tasks compared with the existing baselines.

First, we compare SAPO with three popular gradient-based learning algorithms: CRF, Bi-LSTM, Bi-LSTM-CRF. It is impressive that the proposed SAPO algorithm even has better accuracy than the CRF, Bi-LSTM, and Bi-LSTM-CRF, which are three popular models for sequence labeling. Note that all the models are already fully optimized. As for the superiority, the reason is that the probability is distributed on top- n outputs in SAPO, which is a

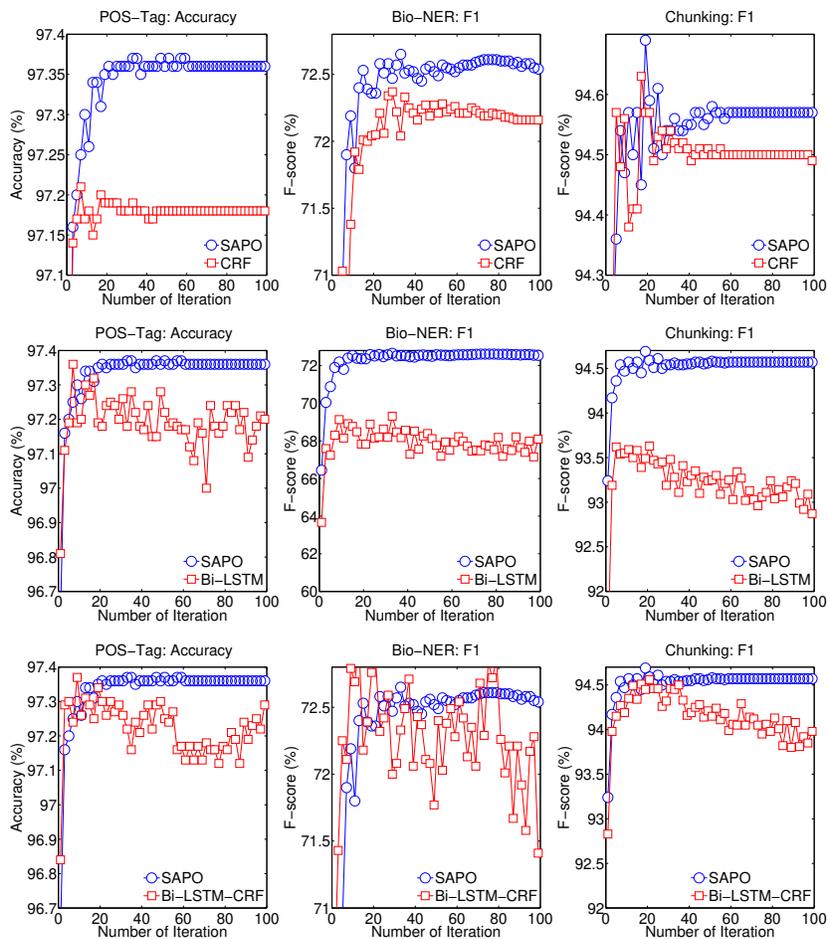


Figure 1: Comparison with CRF, Bi-LSTM, and Bi-LSTM-CRF.

“regularized” distribution instead of the probability distribution spread over all possible outputs (an exponential number). In this sense, SAPO is “regularizing” the exponential probability distribution to a simpler top- n probability distribution. This can be seen as a probability-based regularizer with hyper-parameter n controlling the regularization strength. Interestingly, the experimental results suggest that this type of regularization can indeed improve the accuracy/F-score.

We observe that SAPO is better than CRF, Bi-LSTM, and Bi-LSTM-CRF in all of the three tasks, and it shows that in many cases the differences are

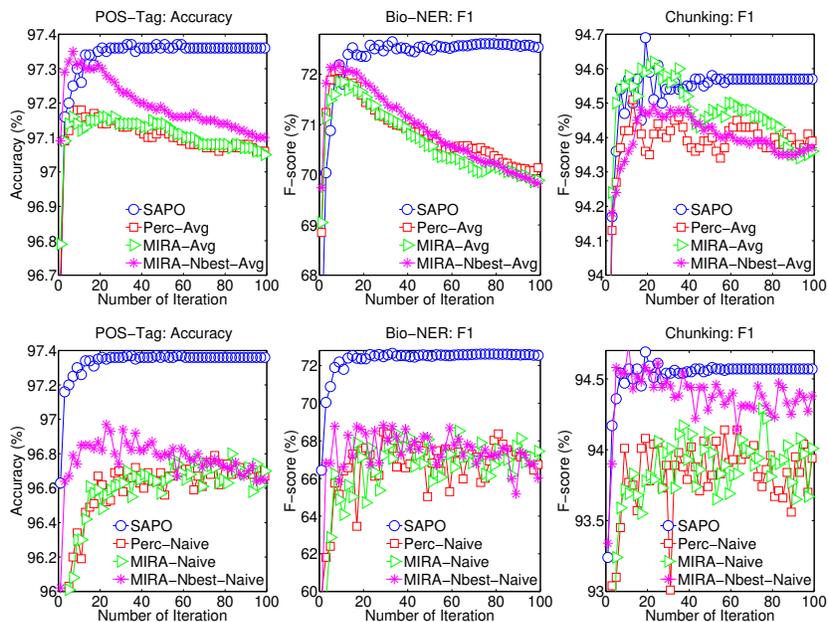


Figure 2: Comparison with structured perceptron and MIRA.

statistically significant. Also, we can see that SAPO is several times faster than CRF and Bi-LSTM in terms of training time. On convergence state, SAPO performs similar or even better than CRF.

Second, we compare SAPO with search-based learning methods, including naive/average versions of Perceptron, MIRA, and Nbest MIRA. As we can see, the superiorities of SAPO over search-based learning methods are even more significant than over CRF.

We also conduct significance tests based on t-test.⁷ For the POS-Tag and chunking task, the significance test suggests that the superiorities of SAPO over all of the baselines except CRF are statistically significant, with at least $p < 0.01$. For the Bio-NER task, the significance test suggests that the superiorities of SAPO over all of the baselines are significant, with at least $p < 0.05$.

Figure 1 and Figure 2 show the training curves based on the number of train-

⁷For the tasks measured by F-score, the t-test is approximated by using accuracy to approximate F-score.

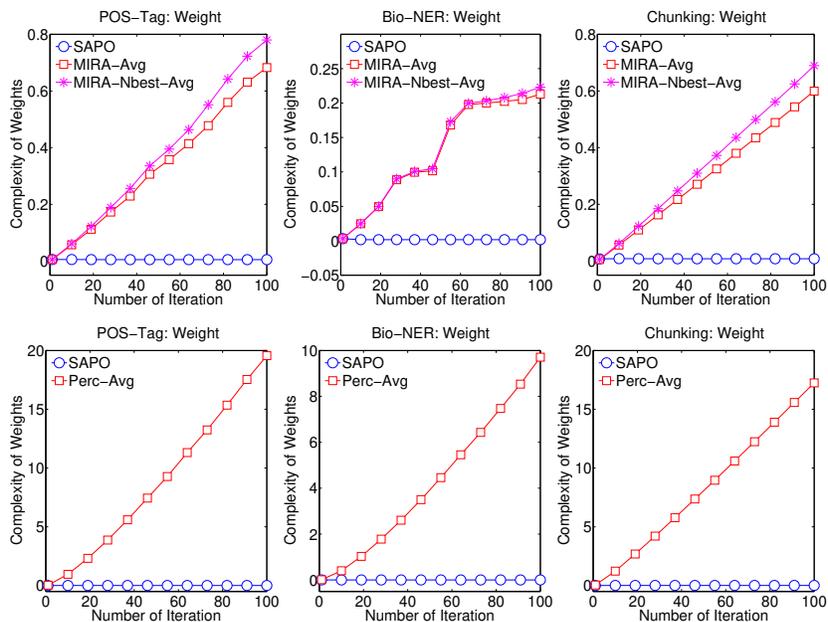


Figure 3: Comparison on parameter weights with structured perceptron and MIRA.

ing iterations. As we can see, SAPO and CRF are convergent as the training goes on, and Bi-LSTM, Bi-LSTM-CRF, Perc, MIRA, and Nbest MIRA diverge as the training goes on.

Figure 3 shows the w -complexity based on the number of training iterations. The w -complexity is the averaged (absolute) value of the weights. As we can see, SAPO is convergent and has very small weight complexity as the training goes on, and Perceptron, MIRA, and Nbest MIRA have linear or even super-linear explosion of weight complexity as the training goes on. Big weight complexity is typically a bad sign for controlling generalization risk.

Figure 4 and Figure 5 show the training time per iteration in terms of seconds. As we can see, SAPO is with low computational cost, especially compared with CRF and Bi-LSTM.

To summarize, the experiment results demonstrate that SAPO has better accuracy than probabilistic gradient-based methods like CRF and BiLSTM, at the same time with fast training speed as structured perceptrons and MIRA.

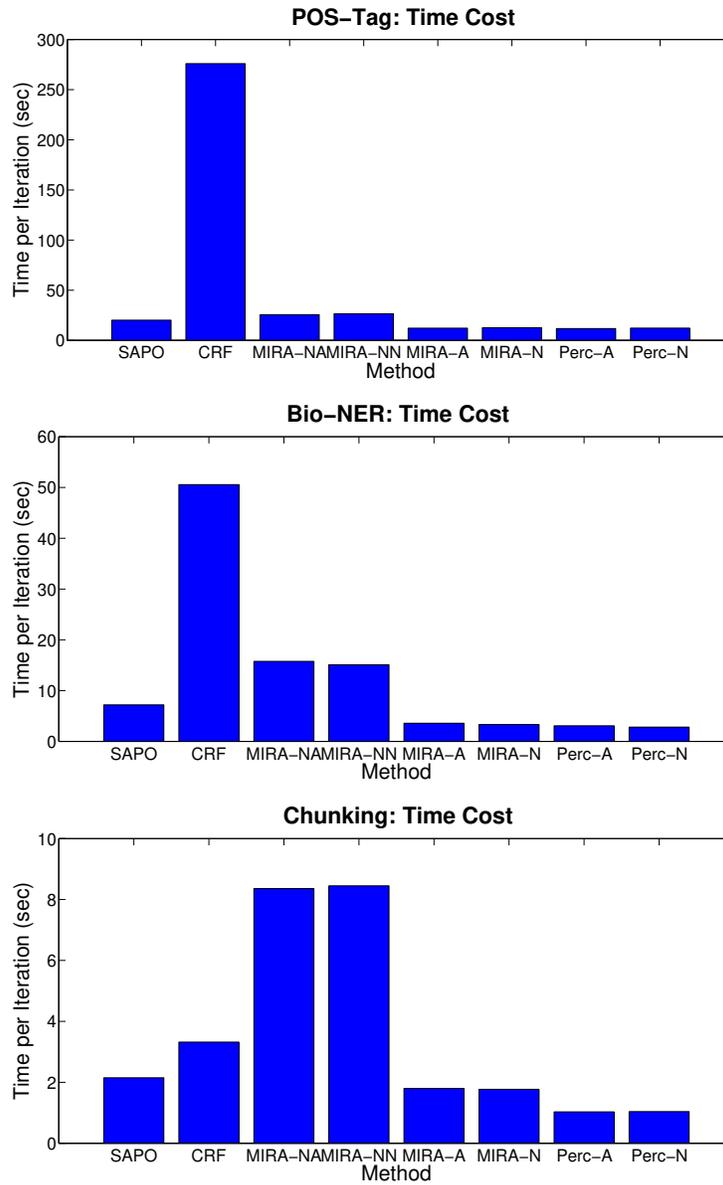


Figure 4: Computation time cost of different models (MIRA-NA: Average Nbest MIRA, MIRA-NN: Naive Nbest MIRA, MIRA-A: Average MIRA, MIRA-N: Naive MIRA, Perc-A: Average Perceptron, Perc-N: Naive Perceptron).

Compared with structured perceptron, which only uses the gold sequence and causes insufficient estimation of the parameters, leading to final divergence,

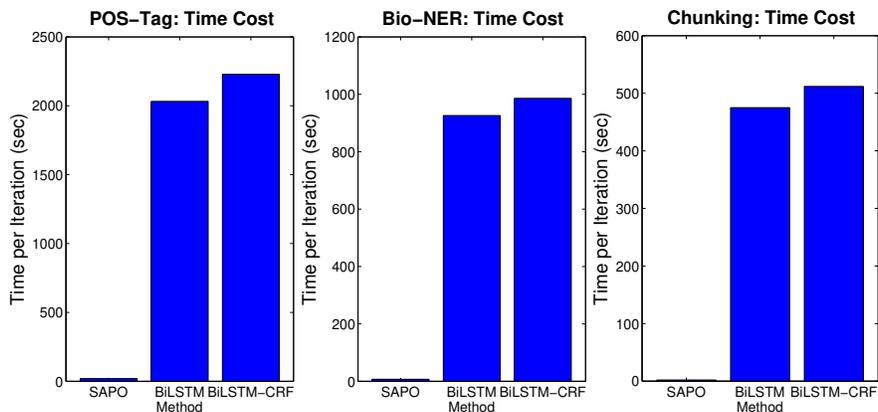


Figure 5: Computation time cost of SAPO, BiLSTM, and BiLSTM-CRF.

SAPO uses top- n possible sequences with probabilistic information in training, and the parameters are more accurate. Compared with CRF, which considers all the possible sequences but causes the overfitting problem, the n in SAPO serves as a kind of regularizer and can control the complexity of the parameters. Compared with LSTM, which is not a structured prediction model at all and does not consider the structural information, SAPO is a structured prediction model and considers the structural information. Moreover, since LSTM can be combined with perceptrons⁸, structured perceptrons and CRF, SAPO can also be combined with LSTM to learn the deep semantic features. Also, SAPO is convergent towards optimum with controllable weight complexity as the training goes on. Note that there are other important advantages of SAPO that are not revealed in those experiments — SAPO supports search-based learning which makes gradient information not necessary and gives probability information, and it is very easy to implement.

5. Related work

Many sequence labeling methods have been developed including probabilistic gradient-based learning methods and search-based learning methods [44, 11].

⁸Most of the LSTM models for sequence labelling belong to this category.

The probabilistic gradient-based learning methods include conditional random fields [16, 38], and a variety of extensions such as dynamic conditional random fields [35], hidden conditional random fields [26], and latent-dynamic conditional random fields [21].

The search-based learning methods include margin infused relaxed algorithm [7], structured perceptrons [5], and a variety of related work in this direction such as latent structured perceptrons [33, 32], confidence weighted linear classification (CW) [9], max-violation perceptrons [45]. Most of the search-based learning methods are large-margin online learning methods. Other related work on sequence labeling also includes maximum margin Markov networks [36] and structured support vector machines [37].

For training sequence labeling models, especially probabilistic gradient-based learning methods like CRF and its variation models, the arguably most popular training method is stochastic gradient descent (SGD) [2, 47, 23, 34, 31], which typically has faster convergence rate compared with alternative batch training methods, such as limited-memory BFGS (L-BFGS) [24] and other quasi-Newton optimization methods. The SGD training has theoretical guarantees to converge to the optimum weights given the convex objective function (e.g., CRF) [2, 47, 23]. For the search-based learning methods such as structured perceptrons, MIRA, and their variation algorithms, the training scheme is usually quite simple and self-contained in the search-based learning algorithm/model [5, 7, 19].

As for dealing with overfitting, the probabilistic gradient-based learning methods typically use explicit regularization terms such as the widely used L_2 regularizer. Other regularization schemes include the L_1 regularizer [1, 39], the group Lasso regularization [46, 18], the structure regularization [30], and others [25]. For the search-based learning methods like structured perceptrons and MIRA, the scheme to deal with overfitting is less formal compared with a regularizer, usually by using parameter averaging or voting [5, 19, 8, 3].

6. Conclusions and future work

The existing sequence labeling methods are problematic. The existing probabilistic gradient-based methods such as CRF and LSTM have slow training speed and do not support search-based optimization. The existing search-based learning methods such as structured perceptrons and MIRA have relatively low accuracy and are non-convergent in most of the real-world tasks. We propose a novel and “easy” solution, a search-based probabilistic online learning framework SAPO, to address most of those issues. SAPO is with fast training, able to support search-based optimization, very easy to implement, with top accuracy, with probabilistic information, and with theoretical guarantees of convergence.

Experiments on well-known benchmark tasks demonstrate that SAPO has better accuracy than CRF and BiLSTM and roughly comparable training speed as structured perceptrons and MIRA. Results also show that SAPO can easily beat the strong baseline systems on those competitive tasks.

In the current implementation, our top- n search uses a simple A* search algorithm with Viterbi heuristics. This top- n search algorithm is not fully optimized for speed. There are several other top- n search algorithms possibly with faster speed. In the future we can optimize the top- n search algorithm. We believe that this can further improve the training speed of SAPO. Moreover, SAPO is a general purpose algorithm for structured classification with arbitrary structures. In the future we can apply SAPO to structured classification with structures that are more complex, e.g., syntactic parsing [17] and statistical machine translation [43].

Acknowledgements

We thank the anonymous reviewers for their thoughtful comments. This work was supported in part by National Natural Science Foundation of China (No. 61673028).

References

References

- [1] Andrew, G., Gao, J., 2007. Scalable training of L_1 -regularized log-linear models, in: International Conference on Machine Learning (ICML), pp. 33–40.
- [2] Bottou, L., LeCun, Y., 2003. Large scale online learning., in: Advances in Neural Information Processing Systems (NIPS), MIT Press.
- [3] Chiang, D., 2012. Hope and fear for discriminative training of statistical translation models. *Journal of Machine Learning Research* 13, 1159–1187.
- [4] Cho, H., Okazaki, N., Miwa, M., Tsujii, J., 2013. Named entity recognition with multiple segment representations. *Inf. Process. Manage.* 49, 954–965.
- [5] Collins, M., 2002. Discriminative training methods for hidden markov models: Theory and experiments with perceptron algorithms, in: Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 1–8.
- [6] Collobert, R., Weston, J., Bottou, L., Karlen, M., Kavukcuoglu, K., Kuksa, P.P., 2011. Natural language processing (almost) from scratch. *Journal of Machine Learning Research* 12, 2493–2537.
- [7] Crammer, K., Singer, Y., 2003. Ultraconservative online algorithms for multiclass problems. *Journal of Machine Learning Research* 3, 951–991.
- [8] Daumé III, H., 2006. Practical Structured Learning Techniques for Natural Language Processing. Ph.D. thesis. University of Southern California.
- [9] Dredze, M., Crammer, K., Pereira, F., 2008. Confidence-weighted linear classification., in: International Conference on Machine Learning (ICML), pp. 264–271.
- [10] Freund, Y., Schapire, R., 1999. Large margin classification using the perceptron algorithm. *Machine Learning* 37, 277–296.

- [11] Fujii, R., Domoto, R., Mochihashi, D., 2017. Nonparametric bayesian semi-supervised word segmentation. *TACL* 5, 179–189.
- [12] Hart, P., Nilsson, N., Raphael, B., 1968. A formal basis for the heuristic determination of minimum cost path. *IEEE Trans. On System Science and Cybernetics* SSC-4(2), 100–107.
- [13] Hochreiter, S., Schmidhuber, J., 1997. Long Short-Term Memory 9, 1735–1780.
- [14] Huang, Z., Xu, W., Yu, K., 2015. Bidirectional LSTM-CRF models for sequence tagging. *arXiv preprint arXiv:1508.01991*.
- [15] Kingma, D.P., Ba, J., 2014. Adam: A method for stochastic optimization. *CoRR* abs/1412.6980.
- [16] Lafferty, J., McCallum, A., Pereira, F., 2001. Conditional random fields: Probabilistic models for segmenting and labeling sequence data, in: *International Conference on Machine Learning (ICML)*, pp. 282–289.
- [17] Li, S., Wang, L., Cao, Z., Li, W., 2014. Text-level discourse dependency parsing, in: *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics, ACL 2014, June 22-27, 2014, Baltimore, MD, USA, Volume 1: Long Papers*, pp. 25–35.
- [18] Martins, A.F.T., Smith, N.A., Figueiredo, M.A.T., Aguiar, P.M.Q., 2011. Structured sparsity in structured prediction., in: *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1500–1511.
- [19] McDonald, R.T., Crammer, K., Pereira, F.C.N., 2005. Online large-margin training of dependency parsers., in: *Annual Meeting of the Association for Computational Linguistics (ACL)*.
- [20] Miwa, M., Thompson, P., McNaught, J., Kell, D.B., Ananiadou, S., 2012. Extracting semantically enriched events from biomedical literature. *BMC Bioinformatics* 13, 108.

- [21] Morency, L.P., Quattoni, A., Darrell, T., 2007. Latent-dynamic discriminative models for continuous gesture recognition, in: Proceedings of CVPR'07, pp. 1–8.
- [22] Nemirovski, A., Juditsky, A., Lan, G., Shapiro, A., 2009. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization* 19, 1574–1609.
- [23] Niu, F., Recht, B., Re, C., Wright, S.J., 2011. Hogwild: A lock-free approach to parallelizing stochastic gradient descent., in: *Advances in Neural Information Processing Systems (NIPS)*, pp. 693–701.
- [24] Nocedal, J., Wright, S.J., 1999. *Numerical optimization*. Springer .
- [25] Quattoni, A., Carreras, X., Collins, M., Darrell, T., 2009. An efficient projection for l_1 regularization., in: *International Conference on Machine Learning (ICML)*, p. 108.
- [26] Quattoni, A., Wang, S., Morency, L.P., Collins, M., Darrell, T., 2007. Hidden conditional random fields. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 29, 1848–1852.
- [27] Sang, E.T.K., Buchholz, S., 2000. Introduction to the CoNLL-2000 shared task: Chunking, in: *Proceedings of CoNLL'00*, pp. 127–132.
- [28] Schuster, M., Paliwal, K.K., 1997. Bidirectional recurrent neural networks. *IEEE Transactions on Signal Processing* 45, 2673–2681.
- [29] Sha, F., Pereira, F., 2003. Shallow parsing with conditional random fields, in: *Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pp. 134–141.
- [30] Sun, X., 2014. Structure regularization for structured prediction, in: *Advances in Neural Information Processing Systems (NIPS)*, pp. 2402–2410.

- [31] Sun, X., Li, W., Wang, H., Lu, Q., 2014. Feature-frequency-adaptive on-line training for fast and accurate natural language processing. *Computational Linguistics* 40, 563–586.
- [32] Sun, X., Matsuzaki, T., Li, W., 2013. Latent structured perceptrons for large-scale learning with hidden information. *IEEE Trans. Knowl. Data Eng.* 25, 2063–2075.
- [33] Sun, X., Matsuzaki, T., Okanohara, D., Tsujii, J., 2009. Latent variable perceptron algorithm for structured classification, in: *International Joint Conference on Artificial Intelligence (IJCAI)*, pp. 1236–1242.
- [34] Sun, X., Wang, H., Li, W., 2012. Fast online training with frequency-adaptive learning rates for chinese word segmentation and new word detection., in: *Annual Meeting of the Association for Computational Linguistics (ACL)*, pp. 253–262.
- [35] Sutton, C., Rohanimanesh, K., McCallum, A., 2004. Dynamic conditional random fields: Factorized probabilistic models for labeling and segmenting sequence data, in: *International Conference on Machine Learning (ICML)*.
- [36] Taskar, B., Guestrin, C., Koller, D., 2003. Max-margin markov networks, in: *Advances in Neural Information Processing Systems (NIPS)*.
- [37] Tsochantaridis, I., Hofmann, T., Joachims, T., Altun, Y., 2004. Support vector machine learning for interdependent and structured output spaces, in: *International Conference on Machine Learning (ICML)*, pp. 823–830.
- [38] Tsuruoka, Y., Tsujii, J., Ananiadou, S., 2009a. Fast full parsing by linear-chain conditional random fields, in: *EACL 2009, 12th Conference of the European Chapter of the Association for Computational Linguistics, Proceedings of the Conference, Athens, Greece, March 30 - April 3, 2009*, pp. 790–798.
- [39] Tsuruoka, Y., Tsujii, J., Ananiadou, S., 2009b. Stochastic gradient descent training for l1-regularized log-linear models with cumulative penalty, in:

Annual Meeting of the Association for Computational Linguistics (ACL), pp. 477–485.

- [40] Uchiumi, K., Tsukahara, H., Mochihashi, D., 2015. Inducing word and part-of-speech with pitman-yor hidden semi-markov models, in: Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing, ACL 2015, July 26-31, 2015, Beijing, China, Volume 1: Long Papers, pp. 1774–1782.
- [41] Vishwanathan, S., Schraudolph, N.N., Schmidt, M.W., Murphy, K.P., 2006. Accelerated training of conditional random fields with stochastic meta-descent, in: International Conference on Machine Learning (ICML), pp. 969–976.
- [42] Viterbi, A.J., 1967. Error bounds for convolutional codes and an asymptotically optimum decoding algorithm. *IEEE Transactions on Information Theory* 13(2), 260–269.
- [43] Wang, R., Zhao, H., Ploux, S., Lu, B., Utiyama, M., 2016. A bilingual graph-based semantic model for statistical machine translation, in: Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence, IJCAI 2016, New York, NY, USA, 9-15 July 2016, pp. 2950–2956.
- [44] Wu, M., Li, W., Lu, Q., Li, B., 2005. CTEMP: A chinese temporal parser for extracting and normalizing temporal information, in: Natural Language Processing - IJCNLP 2005, Second International Joint Conference, Jeju Island, Korea, October 11-13, 2005, Proceedings, pp. 694–706.
- [45] Yu, H., Huang, L., Mi, H., Zhao, K., 2013. Max-violation perceptron and forced decoding for scalable mt training., in: Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 1112–1123.

- [46] Yuan, M., Lin, Y., 2006. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society, Series B* 68, 49–67.
- [47] Zinkevich, M., Weimer, M., Smola, A.J., Li, L., 2010. Parallelized stochastic gradient descent., in: *Advances in Neural Information Processing Systems (NIPS)*, pp. 2595–2603.

A. Appendix: Proof

Here we give the proof of Theorem 1. First, the recursion formula is derived. Then, the bounds are derived.

A.1. Recursion Formula

By subtracting \mathbf{w}^* from both sides and taking norms for (15), we have

$$\begin{aligned} \|\mathbf{w}_{t+1} - \mathbf{w}^*\|^2 &= \|\mathbf{w}_t - \gamma \mathbf{s}_{\mathbf{z}_t}(\mathbf{w}_t) - \mathbf{w}^*\|^2 \\ &= \|\mathbf{w}_t - \mathbf{w}^*\|^2 - 2\gamma(\mathbf{w}_t - \mathbf{w}^*)^T \mathbf{s}_{\mathbf{z}_t}(\mathbf{w}_t) + \gamma^2 \|\mathbf{s}_{\mathbf{z}_t}(\mathbf{w}_t)\|^2 \end{aligned} \quad (25)$$

Taking expectations and let $a_t = \mathbb{E}\|\mathbf{w}_t - \mathbf{w}^*\|^2$, we have

$$\begin{aligned} a_{t+1} &= a_t - 2\gamma \mathbb{E}[(\mathbf{w}_t - \mathbf{w}^*)^T \mathbf{s}_{\mathbf{z}_t}(\mathbf{w}_t)] + \gamma^2 \mathbb{E}[\|\mathbf{s}_{\mathbf{z}_t}(\mathbf{w}_t)\|^2] \\ &\quad \text{(based on (18))} \\ &\leq a_t - 2\gamma \mathbb{E}[(\mathbf{w}_t - \mathbf{w}^*)^T \mathbf{s}_{\mathbf{z}_t}(\mathbf{w}_t)] + \gamma^2 \kappa^2 \\ &\quad \text{(since the random draw of } \mathbf{z}_t \text{ is independent of } \mathbf{w}_t) \\ &= a_t - 2\gamma \mathbb{E}[(\mathbf{w}_t - \mathbf{w}^*)^T \mathbb{E}_{\mathbf{z}_t}(\mathbf{s}_{\mathbf{z}_t}(\mathbf{w}_t))] + \gamma^2 \kappa^2 \\ &= a_t - 2\gamma \mathbb{E}[(\mathbf{w}_t - \mathbf{w}^*)^T \mathbf{s}(\mathbf{w}_t)] + \gamma^2 \kappa^2 \end{aligned} \quad (26)$$

We define

$$\delta(\mathbf{w}) = \nabla f(\mathbf{w}) - \mathbf{s}(\mathbf{w}) \quad (27)$$

and insert it into (16), it goes to

$$\begin{aligned} f(\mathbf{w}') &\geq f(\mathbf{w}) + (\mathbf{w}' - \mathbf{w})^T [\mathbf{s}(\mathbf{w}) + \delta(\mathbf{w})] + \frac{c}{2} \|\mathbf{w}' - \mathbf{w}\|^2 \\ &= f(\mathbf{w}) + (\mathbf{w}' - \mathbf{w})^T \mathbf{s}(\mathbf{w}) + \frac{c}{2} \|\mathbf{w}' - \mathbf{w}\|^2 + (\mathbf{w}' - \mathbf{w})^T \delta(\mathbf{w}) \end{aligned} \quad (28)$$

By setting $\mathbf{w}' = \mathbf{w}^*$, we further have

$$\begin{aligned} (\mathbf{w} - \mathbf{w}^*)^T \mathbf{s}(\mathbf{w}) &\geq f(\mathbf{w}) - f(\mathbf{w}^*) + \frac{c}{2} \|\mathbf{w} - \mathbf{w}^*\|^2 - (\mathbf{w} - \mathbf{w}^*)^T \delta(\mathbf{w}) \\ &\geq \frac{c}{2} \|\mathbf{w} - \mathbf{w}^*\|^2 - (\mathbf{w} - \mathbf{w}^*)^T \delta(\mathbf{w}) \end{aligned} \quad (29)$$

Combining (26) and (29), we have

$$\begin{aligned} a_{t+1} &\leq a_t - 2\gamma \mathbb{E} \left[\frac{c}{2} \|\mathbf{w}_t - \mathbf{w}^*\|^2 - (\mathbf{w}_t - \mathbf{w}^*)^T \delta(\mathbf{w}_t) \right] + \gamma^2 \kappa^2 \\ &= (1 - c\gamma) a_t + 2\gamma \mathbb{E} [(\mathbf{w}_t - \mathbf{w}^*)^T \delta(\mathbf{w}_t)] + \gamma^2 \kappa^2 \end{aligned} \quad (30)$$

Considering (20) and (27), it goes to

$$a_{t+1} \leq (1 - c\gamma) a_t + 2\gamma\tau + \gamma^2 \kappa^2 \quad (31)$$

We can find a steady state a_∞ as follows

$$a_\infty = (1 - c\gamma) a_\infty + 2\gamma\tau + \gamma^2 \kappa^2 \quad (32)$$

which gives

$$a_\infty = \frac{2\tau + \gamma\kappa^2}{c} \quad (33)$$

Defining the function $A(x) = (1 - c\gamma)x + 2\gamma\tau + \gamma^2 \kappa^2$, based on (31) we have

$$\begin{aligned} a_{t+1} &\leq A(a_t) \\ &\quad (\text{Taylor expansion of } A(\cdot) \text{ based on } a_\infty, \text{ with } \nabla^2 A(\cdot) \text{ being } 0) \\ &= A(a_\infty) + \nabla A(a_\infty)(a_t - a_\infty) \\ &= A(a_\infty) + (1 - c\gamma)(a_t - a_\infty) \\ &= a_\infty + (1 - c\gamma)(a_t - a_\infty) \end{aligned} \quad (34)$$

Thus, we have

$$a_{t+1} - a_\infty \leq (1 - c\gamma)(a_t - a_\infty) \quad (35)$$

Unwrapping (35) goes to

$$a_t \leq (1 - c\gamma)^t (a_0 - a_\infty) + a_\infty \quad (36)$$

A.2. Bounds

Since $\nabla f(\mathbf{w})$ is Lipschitz according to (17), we have

$$f(\mathbf{w}) \leq f(\mathbf{w}') + \nabla f(\mathbf{w}')^T(\mathbf{w} - \mathbf{w}') + \frac{q}{2}\|\mathbf{w} - \mathbf{w}'\|^2$$

Setting $\mathbf{w}' = \mathbf{w}^*$, it goes to $f(\mathbf{w}) - f(\mathbf{w}^*) \leq \frac{q}{2}\|\mathbf{w} - \mathbf{w}^*\|^2$, such that

$$\mathbb{E}[f(\mathbf{w}_t) - f(\mathbf{w}^*)] \leq \frac{q}{2}\mathbb{E}\|\mathbf{w}_t - \mathbf{w}^*\|^2 = \frac{q}{2}a_t$$

In order to have

$$E[f(\mathbf{w}_t) - f(\mathbf{w}^*)] \leq \epsilon \tag{37}$$

it is required that $\frac{q}{2}a_t \leq \epsilon$, that is

$$a_t \leq \frac{2\epsilon}{q} \tag{38}$$

Combining (36) and (38), it is required that

$$(1 - c\gamma)^t(a_0 - a_\infty) + a_\infty \leq \frac{2\epsilon}{q} \tag{39}$$

To meet this requirement, it is sufficient to set the learning rate γ such that both terms on the left side are less than $\frac{\epsilon}{q}$. For the requirement of the second term $a_\infty \leq \frac{\epsilon}{q}$, recalling (33), it goes to

$$\gamma \leq \frac{c\epsilon - 2\tau q}{q\kappa^2}$$

Thus, introducing a real value $\beta \geq 1$, we can set γ as

$$\gamma = \frac{c\epsilon - 2\tau q}{\beta q\kappa^2} \tag{40}$$

Note that, to make this formula meaningful, it is required that

$$c\epsilon - 2\tau q \geq 0$$

Thus, it is required that

$$\tau \leq \frac{c\epsilon}{2q}$$

which is solved by the condition of (21).

On the other hand, we analyze the requirement of the first term that

$$(1 - c\gamma)^t (a_0 - a_\infty) \leq \frac{\epsilon}{q} \quad (41)$$

Since $a_0 - a_\infty \leq a_0$, it holds by requiring

$$(1 - c\gamma)^t a_0 \leq \frac{\epsilon}{q} \quad (42)$$

which goes to

$$t \geq \frac{\log \frac{\epsilon}{qa_0}}{\log(1 - c\gamma)} \quad (43)$$

Since $\log(1 - c\gamma) \leq -c\gamma$ given (19), and that $\log \frac{\epsilon}{qa_0}$ is a negative term, we have

$$\frac{\log \frac{\epsilon}{qa_0}}{\log(1 - c\gamma)} \leq \frac{\log \frac{\epsilon}{qa_0}}{-c\gamma}$$

Thus, (43) holds by requiring

$$\begin{aligned} t &\geq \frac{\log \frac{\epsilon}{qa_0}}{-c\gamma} \\ &= \frac{\log(qa_0/\epsilon)}{c\gamma} \end{aligned} \quad (44)$$

Combining (40) and (44), it goes to

$$t \geq \frac{\beta q \kappa^2 \log(qa_0/\epsilon)}{c(c\epsilon - 2\tau q)}$$

which completes the proof.