

# Search-based structured prediction

Hal Daumé III · John Langford · Daniel Marcu

Received: 22 September 2006 / Revised: 15 May 2008 / Accepted: 16 January 2009 /  
Published online: 14 March 2009  
Springer Science+Business Media, LLC 2009

**Abstract** We present SEARN, an algorithm for integrating SEARCH and LEARNING to solve complex structured prediction problems such as those that occur in natural language, speech, computational biology, and vision. SEARN is a meta-algorithm that transforms these complex problems into simple classification problems to which any binary classifier may be applied. Unlike current algorithms for structured learning that require *decomposition* of both the loss function and the feature functions over the predicted structure, SEARN is able to learn prediction functions for *any* loss function and *any* class of features. Moreover, SEARN comes with a strong, natural theoretical guarantee: good performance on the derived classification problems implies good performance on the structured prediction problem.

**Keywords** Structured prediction · Search · Reductions

## 1 Introduction

Prediction is the task of learning a function  $f$  that maps inputs  $x$  in an input domain  $\mathcal{X}$  to outputs  $y$  in an output domain  $\mathcal{Y}$ . Standard algorithms—support vector machines, decision trees, neural networks, etc.—focus on “simple” output domains such as  $\mathcal{Y} = \{-1, +1\}$  (in the case of binary classification) or  $\mathcal{Y} = \mathbb{R}$  (in the case of univariate regression).

We are interested in problems for which elements  $y \in \mathcal{Y}$  have complex internal structure. The simplest and best studied such output domain is that of labeled sequences. However, we are interested in even more complex domains, such as the space of English sentences

---

Editor: Dan Roth.

H. Daumé III (✉)  
School of Computing, University of Utah, Salt Lake City, UT 84112, USA  
e-mail: [me@hal3.name](mailto:me@hal3.name)

J. Langford  
Yahoo! Research Labs, New York, NY 10011, USA

D. Marcu  
Information Sciences Institute, Marina del Rey, CA 90292, USA

(for instance in a machine translation application), the space of short documents (perhaps in an automatic document summarization application), or the space of possible assignments of elements in a database (in an information extraction/data mining application). The structured complexity of features and loss functions in these problems significantly exceeds that of sequence labeling problems.

From a high level, there are four dimensions along which structured prediction algorithms vary: structure (varieties of structure for which efficient learning is possible), loss (different loss functions for which learning is possible), features (generality of feature functions for which learning is possible) and data (ability of algorithm to cope with imperfect data sources such as missing data, etc.). An in-depth discussion of alternative structured prediction algorithms is given in Sect. 5. However, to give a flavor, the popular conditional random field algorithm (Lafferty et al. 2001) is viewed along these dimensions as follows. Structure: inference for a CRF is tractable for any graphical model with bounded tree width; Loss: the CRF typically optimizes a log-loss approximation to 0/1 loss over the entire structure; Features: any feature of the input is possible but only output features that obey the graphical model structure are allowed; Data: EM can cope with hidden variables.

We prefer a structured prediction algorithm that is not limited to models with bounded treewidth, is applicable to any loss function, can handle arbitrary features and can cope with imperfect data. Somewhat surprisingly, SEARN meets nearly all of these requirements by transforming structured prediction problems into binary prediction problems to which a vanilla binary classifier can be applied. SEARN comes with a strong theoretical guarantee: good binary classification performance implies good structured prediction performance. Simple applications of SEARN to standard structured prediction problems yield tractable state-of-the-art performance. Moreover, we can apply SEARN to more complex, non-standard structured prediction problems and achieve excellent empirical performance.

This paper has the following outline:

1. Introduction.
2. Core Definitions.
3. The SEARN Algorithm.
4. Theoretical Analysis.
5. A Comparison to Alternative Techniques.
6. Experimental results.
7. Discussion.

## 2 Core definitions

In order to proceed, it is useful to formally define a structured prediction problem in terms of a state space.

**Definition 1** A *structured prediction* problem  $\mathcal{D}$  is a cost-sensitive classification problem where  $\mathcal{Y}$  has structure: elements  $y \in \mathcal{Y}$  decompose into variable-length vectors  $(y_1, y_2, \dots, y_T)$ .<sup>1</sup>  $\mathcal{D}$  is a distribution over inputs  $x \in \mathcal{X}$  and cost vectors  $\mathbf{c}$ , where  $|\mathbf{c}|$  is a variable in  $2^T$ .

---

<sup>1</sup>Treating  $y$  as a vector is simply a useful encoding; we are not interested only in sequence labeling problems.

As a simple example, consider a parsing problem under  $F_1$  (balanced precision/recall) loss. In this case,  $\mathcal{D}$  is a distribution over  $(x, c)$  where  $x$  is an input sequence and for all trees  $y$  with  $|x|$ -many leaves,  $c_y$  is the  $F_1$  loss of  $y$  against the “true” output.

The goal of structured prediction is to find a function  $h : \mathcal{X} \rightarrow \mathcal{Y}$  that minimizes the loss given in (1):

$$L(\mathcal{D}, h) = \mathbb{E}_{(x,c) \sim \mathcal{D}} \{c_{h(x)}\}. \quad (1)$$

The algorithm we present is based on the view that a vector  $y \in \mathcal{Y}$  can be produced by predicting each component  $(y_1, \dots, y_T)$  in turn, allowing for dependent predictions. This is important for coping with general loss functions.

For a data set  $(x_1, c_1), \dots, (x_N, c_N)$  of structured prediction examples, we write  $T_n$  for the length of the longest search path on example  $n$ , and  $T_{\max} = \max_n T_n$ .

### 3 The SEARN algorithm

The core idea behind SEARN is as follows. We are given a set of structured prediction examples of the form  $(x, c)$ . We view the production of a structured output as a *search* process over the decomposition of  $y$  into  $y_1, \dots, y_T$ . We train a (cost-sensitive) classifier to predict each of the  $y_i$  components in turn; this is akin to a greedy search process (an extension to non-greedy search is considered in Sect. 3.4.3). The  $t$ th decision may be dependent on any of the preceding  $t - 1$  decisions. The cost-sensitive classifier that makes each of these search decisions should be trained to do well (with respect to  $c$ ) for any prefix of  $t - 1$  decisions that it might encounter. This introduces a chicken-and-egg problem: how to best train the classifier depends on the classifier itself. We solve this using an iterative scheme.

We formulate the SEARN algorithm in terms of a *policy*, a notion borrowed from reinforcement learning. A policy tells us: for a given “node” in the search space, what is the best action to take? It is important to keep in mind, however, the separation between *training* data and *test* data in the structured prediction setting. In particular, what we desire is a good policy *for the test data*. We assume that we already have a good policy for the training data: this policy can be constructed on the basis of the known loss vector  $c$  for any particular training example.

For example, consider a simple sequence labeling task (for instance, part-of-speech tagging on natural language sentences). Here, the output is a vector of labels and for simplicity we assume that the loss is Hamming loss (number of incorrectly predicted labels). Given a test example, we wish to predict each label in the sequence incrementally. For instance, we may tag the sequence left-to-right, labeling the  $t$  word only after the previous  $t - 1$  words have been labeled. In this case, we associate a node in a search space to be any labeling of a prefix of the sentence. An action in this space corresponds to labeling one additional word. The learned policy predict, for a given node (i.e., prefix of labels) and the sentence itself, what is the best action to take. That is, how should we label the next word. In the case of Hamming loss over sequences, it is trivial to construct a good policy when the correct label sequence is known: the best thing to do is to always predict the correct tag. The *learning* problem is, essentially, to transfer this policy that works when we observe the true label sequence to examples where we do not observe the true label sequence.

More formally, there are several vital ingredients in any application of SEARN: a search space for decomposing the prediction problem; a cost sensitive learning algorithm; labeled structured prediction training data; a known loss function for the structured prediction problem; and a good initial policy. These aspects are described in more detail below.

A search space  $\mathcal{S}$ . The choice of search space plays a role similar to the choice of structured decomposition in other algorithms. Final elements of the search space can always be referenced by a sequence of choices  $\hat{y}$ . In simple applications of SEARN the search space is concrete. For example, it might consist of the parts of speech of each individual word in a sentence. In general, the search space can be abstract, and we show this can be beneficial experimentally (Sect. 3.4.3). An abstract search space comes with an (unlearned) function  $f(\hat{y})$  which turns any sequence of predictions in the abstract search space into an output of the correct form. (For a concrete search space,  $f$  is just the identity function. To minimize confusion, we will leave off  $f$  in future notation unless its presence is specifically important.) We discuss the effect of the search space on experimental performance in Sect. 6.1.8.

A cost sensitive learning algorithm  $A$ . The learning algorithm returns a multiclass classifier  $h(s)$  given cost sensitive training data. Here  $s$  is a description of the location in the search space. A reduction of cost sensitive classification to binary classification (Beygelzimer et al. 2005) reduces the requirement to an importance-weighted binary learning algorithm. SEARN relies upon this learning algorithm to form good generalizations. Nothing else in the SEARN algorithm attempts to achieve generalization or estimation. The performance of SEARN is strongly dependent upon how capable the learned classifier is. We call the learned classifier a *policy* because it is used multiple times on inputs which it effects, just as in reinforcement learning.

Labeled structured prediction training data. SEARN digests the labeled training data for the structured prediction problem into cost-sensitive training data which is fed into the cost-sensitive learning algorithm.<sup>2</sup>

A known loss function. A loss function  $L(y, f(\hat{y}))$  must be known and be computable for any sequence of predictions.

A good initial policy. This policy should achieve low loss when applied to the training data. This can (but need not always) be defined using a search algorithm.

### 3.1 SEARN at test time

SEARN at test time is a very simple algorithm. It uses the policy returned by the learning algorithm to construct a sequence of decisions  $\hat{y}$  and makes a final prediction  $f(\hat{y})$ . First, one uses the learned policy to compute  $y_0$  on the basis of just the input  $x$ . One then computes  $y_1$  on the basis of  $x$  and  $y_0$ , followed by predicting  $y_2$  on the basis of  $x$ ,  $y_0$  and  $y_1$ , etc. Finally, one predicts  $y_T$  on the basis of the input  $x$  and *all* previous decisions.

### 3.2 SEARN at train time

SEARN operates in an iterative fashion. At each iteration it uses a known policy to create new cost-sensitive classification examples. These examples are essentially the classification decisions that a policy would need to get right in order to perform search well. These are used to learn a new classifier, which is interpreted as a new policy. This new policy is interpolated with the old policy and the process repeats.

<sup>2</sup>A  $k$ -class cost-sensitive example is given by an input  $X$  and a vector of costs  $\mathbf{c} \in (\mathbb{R}^+)^k$ . Each class  $i$  has an associated cost  $c_i$  and the goal is a function  $h : X \mapsto i$  that minimizes the expected value of  $c_i$ . See (Beygelzimer et al. 2005).

### 3.2.1 Initial policy

SEARN relies on a good initial policy on the *training data*. This policy can take full advantage of the training data labels. The initial policy needs to be efficiently computable for SEARN to be efficient. The implications of this assumption are discussed in detail in Sect. 3.4.1, but it is strictly weaker than assumptions made by other structured prediction techniques. The initial policy we use is a policy that, for a given state predicts the best action to take with respect to the labels:

**Definition 2** (Initial Policy) For an input  $x$  and a cost vector  $c$  as in Def 1, and a state  $s = x \times (y_1, \dots, y_t)$  in the search space, the *initial policy*  $\pi(s, c)$  is  $\arg \min_{y_{t+1}} \min_{y_{t+2}, \dots, y_T} C_{(y_1, \dots, y_T)}$ . That is,  $\pi$  chooses the action (i.e., value for  $y_{t+1}$ ) that minimizes the corresponding cost, assuming that all *future* decisions are also made optimally.

This choice of initial policy is optimal when the correct output is a deterministic function of the input features (effectively in a noise-free environment).

### 3.2.2 Cost-sensitive examples

In the training phase, SEARN uses a given policy  $h$  (initialized to the initial policy  $\pi$ ) to construct cost-sensitive multiclass classification examples from which a new classifier is learned. These classification examples are created by *running* the given policy  $h$  over the training data. This generates one path per structured training example. SEARN creates a single cost-sensitive example for each state on each path. The classes associated with each example are the available actions (or next states). The only difficulty lies in specifying the costs.

The cost associated with taking an action that leads to state  $s$  is the *regret* associated with this action, given our current policy. For each state  $s$  and each action  $a$ , we take action  $a$  and then execute the policy to gain a full sequence of predictions  $\hat{y}$  for which we can compute a loss  $c_{\hat{y}}$ . Of all the possible actions, one,  $a'$ , has the minimum expected loss. The cost  $\ell_h(c, s, a)$  for an action  $a$  in state  $s$  is the difference in loss between taking action  $a$  and taking the action  $a'$ ; see (2):

$$\ell_h(c, s, a) = \mathbb{E}_{\hat{y} \sim (s, a, h)} c_{\hat{y}} - \min_{a'} \mathbb{E}_{\hat{y} \sim (s, a', h)} c_{\hat{y}}. \quad (2)$$

One complication arises because the policy used may be stochastic. First, the base classifier may be stochastic. Second, the *interpolation* is stochastic. In particular, the method of interpolation is to flip a weighted coin. On heads, the initial policy is called; on tails, the learned policy. There are (at least) three possible ways to deal with randomness.

1. Monte-Carlo sampling: one draws many paths according to  $h$  beginning at  $s'$  and average over the costs.
2. Single Monte-Carlo sampling: draw a single path and use the corresponding cost, with tied randomization as per Pegasus (Ng and Jordan 2000).
3. Approximation: it is often possible to efficiently compute the loss associated with following the initial policy from a given state; when  $h$  is sufficiently good, this may serve as a useful and fast approximation. (This is also the approach described by Langford and Zadrozny 2005.)

The quality of the learned solution depends on the quality of the approximation of the loss. Obtaining Monte-Carlo samples is likely the best solution, but in many cases the approximation is sufficient. An empirical comparison of these options is performed in (Daumé III 2006). Here it is observed that for easy problems (one for which low loss is possible), the approximation performs approximately as well as the alternatives. Moreover, typically the approximation outperforms the single sample approach, likely due to the noise induced by following a single sample.

### 3.2.3 Algorithm

The SEARN algorithm is shown in Fig. 1. As input, the algorithm takes a structured learning data set, an initial policy  $\pi$  and a multiclass cost sensitive learner  $A$ . SEARN operates iteratively, maintaining a current policy hypothesis  $h$  at each iteration. This hypothesis is initialized to the initial policy (step 1).

The algorithm then loops for a number of iterations. In each iteration, it creates a (multi-)set of cost-sensitive examples,  $S$ . These are created by looping over each structured example (step 4). For each example (step 5), the *current policy*  $h$  is used to produce a full output, represented as a sequence of predictions  $y_1, \dots, y_{T_n}$ . From this, states are derived and used to create a single cost-sensitive example (steps 6–14) at each timestep. (We evaluate the role of the iterations in Sect. 6.1.9.)

The first task in creating a cost-sensitive example is to compute the associated feature vector, performed in step 7. This feature vector is based on the current state  $s_t$  which includes the features  $x$  (the creation of the feature vectors is discussed in more detail in Sect. 3.3). The cost vector contains one entry for every possible action  $a$  that can be executed from state  $s_t$ . For each action  $a$ , we compute the expected loss associated with the state  $s_t \oplus a$ : the state arrived at assuming we take action  $a$  (step 10).

#### Algorithm SEARN( $S^{\text{SP}}, \pi, A$ )

```

1: Initialize policy  $h \leftarrow \pi$ 
2: while  $h$  has a significant dependence on  $\pi$  (see Lemma 2) do
3:   Initialize the set of cost-sensitive examples  $S \leftarrow \emptyset$ 
4:   for  $(x, y) \in S^{\text{SP}}$  do
5:     Compute predictions under the current policy  $\hat{y} \sim x, h$ 
6:     for  $t = 1 \dots T_x$  do
7:       Compute features  $\Phi = \Phi(s_t)$  for state  $s_t = (x, y_1, \dots, y_t)$ 
8:       Initialize a cost vector  $c = \langle \rangle$ 
9:       for each possible action  $a$  do
10:        Let the cost  $\ell_a$  for example  $x, c$  at state  $s$  be  $\ell_h(c, s, a)$ 
11:       end for
12:       Add cost-sensitive example  $(\Phi, \ell)$  to  $S$ 
13:     end for
14:   end for
15:   Learn a classifier on  $S$ :  $h' \leftarrow A(S)$ 
16:   Interpolate:  $h \leftarrow \beta h' + (1 - \beta)h$ 
17: end while
18: return  $h_{\text{last}}$  without  $\pi$ 

```

**Fig. 1** Complete SEARN algorithm

SEARN creates a large set of cost-sensitive examples  $S$ . These are fed into any cost-sensitive classification algorithm,  $A$ , to produce a new classifier  $h'$  (step 15). In step 16, SEARN *combines* the newly learned classifier  $h'$  with the current classifier  $h$  to produce a new classifier. This combination is performed through stochastic interpolation with interpolation parameter  $\beta$  (see Sect. 4 for details). The meaning of stochastic interpolation here is: “every time  $h$  is evaluated, a new random number is drawn. If the random number is less than  $\beta$  then  $h'$  is used and otherwise the old  $h$  is used.” SEARN returns the final policy with  $\pi$  removed (step 18) and the stochastic interpolation renormalized.

### 3.3 Feature computations

In step 7 of the SEARN algorithm (Fig. 1), one is required to compute a feature vector  $\Phi$  on the basis of the give state  $s_t$ . In theory, this step is arbitrary. However, the performance of the underlying classification algorithm (and hence the induced structured prediction algorithm) hinges on a good choice for these features. The feature vector  $\Phi(s_t)$  may depend on *any* aspect of the input  $x$  and any past decision. In particular, there is no limitation to a “Markov” dependence on previous decisions.

For concreteness, consider the part-of-speech tagging task: for each word in a sentence, we must assign a single part of speech (e.g., Det, Noun, Verb, etc.). Given a state  $s_t = \langle x, y_1, \dots, y_t \rangle$ , one might compute a sparse feature vector  $\Phi(s_t)$  with zeros everywhere except at positions corresponding to “interesting” aspects of the input. For instance, a feature corresponding to the identity of the  $t + 1$ st word in the sentence would likely be very important (since this is the word to be tagged). Furthermore, a feature corresponding to the value  $y_t$  would likely be important, since we believe that subsequent tags are not independent of previous tags. These features would serve as the input to the cost-sensitive learning algorithm, which would attempt to predict the correct label for the  $t + 1$ st word. This usually corresponds to learning a single weight vector for each class (in a one-versus-all setting) or to learning a single weight vector for each pair of classes (for all-pairs).

### 3.4 Policies

SEARN functions in terms of policies, a notion borrowed from the field of reinforcement learning. This section discusses the nature of the initial policy assumption and the connections to reinforcement learning.

#### 3.4.1 Computability of the initial policy

SEARN relies upon the ability to start with a good initial policy  $\pi$ , defined formally in Definition 2. For many simple problems under standard loss functions, it is straightforward to compute a good policy  $\pi$  in constant time. For instance, consider the sequence labeling problem (discussed further in Sect. 6.1). A standard loss function used in this task is Hamming loss: of all possible positions, how many does our model predict incorrectly. If one performs search left-to-right, labeling one element at a time (i.e., each element of the  $y$  vector corresponds exactly to one label), then  $\pi$  is trivial to compute. Given the correct label sequence,  $\pi$  simply chooses at position  $i$  the correct label at position  $i$ . However, SEARN is not limited to simple Hamming loss. A more complex loss function often considered for the sequence segmentation task is F-score over (correctly labeled) segments. As discussed in Sect. 6.1.3, it is just as easy to compute a good initial policy for this loss function. This is not possible in many other frameworks, due to the non-additivity of F-score. This is independent of the features.

This result—that SEARN can learn under strictly more complex structures and loss functions than other techniques—is not limited to sequence labeling, as demonstrated below in Theorem 1. In order to prove this, we need to formalize what we consider as “other techniques.” We use the max-margin Markov network ( $M^3N$ ) formalism (Taskar et al. 2005) for comparison, since this currently appears to be the most powerful generic framework. In particular, learning in  $M^3N$ s is often tractable for problems that would be #P-hard for conditional random fields. The  $M^3N$  has several components, one of which is the ability to compute a loss-augmented minimization (Taskar et al. 2005). This requirement states that (3) is computable for any input  $x$ , output set  $\mathcal{Y}_x$ , true output  $y$  and weight vector  $\mathbf{w}$ :

$$\text{opt}(\mathcal{Y}_x, y, \mathbf{w}) = \arg \max_{\hat{y} \in \mathcal{Y}_x} \mathbf{w}^\top \Phi(x, \hat{y}) - l(y, \hat{y}). \quad (3)$$

In (3),  $\Phi(\cdot)$  produces a vector of features,  $\mathbf{w}$  is a weight vector and  $l(y, \hat{y})$  is the loss for prediction  $\hat{y}$  when the correct output is  $y$ .

**Theorem 1** *Suppose (3) is computable in time  $T(x)$ ; then the optimal policy is computable in time  $\mathcal{O}(T(x))$ .*

*Proof* We use a vector encoding of  $y$  that maintains the decomposition over the regions used by the  $M^3N$ . Given a prefix  $y_1, \dots, y_t$ , solve *opt* on the future choices (i.e., remove the part of the structure corresponding to the first  $t$  outputs) and removing the “loss” term, which gives us an optimal policy.  $\square$

This theorem shows that any problem solvable by a max margin Markov network is also amenable to a solution via SEARN. One key advantage to SEARN, however, is that it never needs to solve such an “arg max” problem. For example, for complex, loopy structures (such as those found in image segmentation problems or complex natural language processing problems), the “arg max” computation is known to be intractable. However, one may *still* apply SEARN to such problems by ordering the vertices and running greedy search.

### 3.4.2 Search-based policies

The SEARN algorithm and the theory to be presented in Sect. 4 do not require that the initial policy be optimal. SEARN can train against *any* policy. One artifact of this observation is that we can use *search* to create the initial policy.

At any step of SEARN, we need to be able to compute the best next action. That is, given a node in the search space, and the cost vector  $\mathbf{c}$ , we need to compute the best step to take. This is *exactly* the standard search problem: given a node in a search space, we find the shortest path to a goal. By taking the first step along this shortest path, we obtain a good initial policy (assuming this shortest path is, indeed, shortest). This means that when SEARN asks for the best next step, one can execute *any* standard search algorithm to compute this, for cases where a good initial policy is not available analytically.

Given this observation, the requirements of SEARN are reduced: instead of requiring a good initial policy, we simply require that one can perform efficient approximate search.



### 3.4.3 Beyond greedy search

We have presented SEARN as an algorithm that mimics the operations of a *greedy* search algorithm. Real-world experience has shown that often greedy search is insufficient and more complex search algorithms are required. This observation is consistent with the standard view of search (trying to find a shortest path), but nebulous when considered in the context of SEARN. Nevertheless, it is often desirable to allow a model to trade past decisions off future decisions, and this is precisely the purpose of instituting more complex search algorithms.

It turns out that any (non-greedy) search algorithm operating in a search space  $\mathcal{S}$  can be equivalently viewed as a greedy search algorithm operating in an abstract space  $\mathcal{S}^*$  (where the structure of the abstract space is dependent on the original search algorithm). In a general search algorithm (Russell and Norvig 1995), one maintains a *queue* of active states and expands a single state in each search step. After expansion, each resulting child state is enqueued. The ordering (and, perhaps, maximal size) of the queue is determined by the specific search algorithm.

In order to simulate this more complex algorithm as greedy search, we construct the abstract space  $\mathcal{S}^*$  as follows. Each node  $s \in \mathcal{S}^*$  represents a state of the queue. A transition exists between  $s$  and  $s'$  in  $\mathcal{S}^*$  exactly when a particular expansion of an  $\mathcal{S}$ -node in the  $s$ -queue results in the queue becoming  $s'$ . Finally, for each goal state  $g \in \mathcal{S}$ , we augment  $\mathcal{S}^*$  with a single unique goal state  $g^*$ . We insert transitions from  $s \in \mathcal{S}^*$  to  $g^*$  exactly when  $g^* \in s$ . Thus, in order to complete the search process, a goal node must be in the queue and the search algorithm must select this single node.

In general, SEARN makes no assumptions about how the search process is structured. A different search process leads to a different bias in the learning algorithm. It is up to the designer to construct a search process so that (a) a good bias is exhibited and (b) computing a good initial policy is easy. For instance, for some combinatorial problems such as matchings or tours, it is known that left-to-right beam search tends to perform poorly. For these problems, a local hill-climbing search is likely to be more effective since we expect it to render the underlying classification problem simpler.

## 4 Theoretical analysis

SEARN functions by slowly moving *away* from the initial policy (which is available only for the training data) *toward* a fully learned policy. Each iteration of SEARN *degrades* the current policy.<sup>3</sup> The main theorem states that the learned policy is not much worse than the starting (optimal) policy plus a term related to the average cost sensitive loss of the learned classifiers and another term related to the maximum cost sensitive loss. To simplify notation, we write  $T$  for  $T_{\max}$ .

It is important in the analysis to refer explicitly to the error of the classifiers learned during SEARN process. Let  $\text{SEARN}(\mathcal{D}, h)$  denote the distribution over classification problems generated by running SEARN with policy  $h$  on distribution  $\mathcal{D}$ . Also let  $\ell_h^{\text{CS}}(h')$  denote the loss of classifier  $h'$  on the distribution  $\text{SEARN}(\mathcal{D}, h)$ . Let the average cost sensitive loss over

<sup>3</sup>In empirical practice, performance first degrades and then often improves. This positive contribution can be understood in the framework of conservative policy iteration (Kakade and Langford 2002).

$I$  iterations be:

$$\ell_{\text{avg}} = \frac{1}{I} \sum_{i=1}^I \ell_{h_i}^{\text{CS}}(h'_i) \tag{4}$$

where  $h_i$  is the  $i$ th policy and  $h'_i$  is the classifier learned on the  $i$ th iteration.

**Theorem 2** For all  $\mathcal{D}$  with  $c_{\text{max}} = \mathbb{E}_{(x,c) \sim \mathcal{D}} \max_y c_y$  (with  $(x, c)$  as in Definition 1), for all learned cost sensitive classifiers  $h'$ , SEARN with  $\beta = 1/T^3$  and  $2T^3 \ln T$  iterations, outputs a learned policy with loss bounded by:

$$L(\mathcal{D}, h_{\text{last}}) \leq L(\mathcal{D}, \pi) + 2T \ell_{\text{avg}} \ln T + (1 + \ln T) c_{\text{max}}/T$$

The dependence on  $T$  in the second term is due to the cost sensitive loss being an average over  $T$  timesteps while the total loss is a sum. The  $\ln T$  factor is not essential and can be removed using other approaches (Bagnell et al. 2003; Langford and Zadrozny 2005). The advantage of the theorem here is that it applies to an algorithm that naturally copes with variable length  $T$  and yields a smaller amount of computation in practice.

The choices of  $\beta$  and the number of iterations are pessimistic in practice. Empirically, we use a development set to perform a line search minimization to find per-iteration values for  $\beta$  and to decide when to stop iterating. The analytical choice of  $\beta$  is made to ensure that the probability that the newly created policy only makes *one* different choice from the previous policy for any given example is sufficiently low. The choice of  $\beta$  assumes the worst: the newly learned classifier *always* disagrees with the previous policy. In practice, this rarely happens. After the first iteration, the learned policy is typically quite good and only rarely differs from the initial policy. So choosing such a small value for  $\beta$  is unnecessary: even with a higher value, the current classifier often agrees with the previous policy.

The theorem makes clear that the performance of the initial policy is of great interest in controlling final performance. For learning problems where noise is not essential, simply using a policy on the training data which agrees with the labels is an optimal starting policy. In situations where noise is essential, constructing a good initial policy can be difficult, even given the labels available on the training data. The theory here does not specify how this initial policy with good performance is formed—all that it does is show that the final learned policy competes well with whatever initial policy is used.

The proof rests on the following lemma.

**Lemma 1** (Policy degradation) Given a policy  $h$  with loss  $L(\mathcal{D}, h)$ , apply a single iteration of SEARN to learn a classifier  $h'$  with cost-sensitive loss  $\ell_h^{\text{CS}}(h')$ . Create a new policy  $h^{\text{new}}$  by acting according to  $h'$  with probability  $\beta \in (0, 1/T)$  and otherwise acting according to  $h$  at each step. Then, for all  $\mathcal{D}$ , with  $c_{\text{max}} = \mathbb{E}_{(x,c) \sim \mathcal{D}} \max_y c_y$  (with  $(x, c)$  as in Definition 1):

$$L(\mathcal{D}, h^{\text{new}}) \leq L(\mathcal{D}, h) + T\beta \ell_h^{\text{CS}}(h') + \frac{1}{2}\beta^2 T^2 c_{\text{max}}. \tag{5}$$

*Proof* The proof largely follows the proofs of Lemma 6.1 and Theorem 4.1 for conservative policy iteration (Kakade and Langford 2002). The three differences are that (1) we must deal with the finite horizon case; (2) we move *away from* rather than *toward* a good policy; and (3) we expand to higher order.

The proof works by separating three cases depending on whether  $h'$  or  $h$  is called in the process of running  $h^{\text{new}}$ . The easiest case is when  $h'$  is never called. The second case is

when it is called exactly once. The final case is when it is called more than once. Denote these three events by  $c = 0$ ,  $c = 1$  and  $c \geq 2$ , respectively:

$$L(\mathcal{D}, h^{\text{new}}) = Pr(c = 0)L(\mathcal{D}, h^{\text{new}} \| c = 0) + Pr(c = 1)L(\mathcal{D}, h^{\text{new}} \| c = 1) + Pr(c \geq 2)L(\mathcal{D}, h^{\text{new}} \| c \geq 2) \tag{6}$$

$$\leq (1 - \beta)^T L(\mathcal{D}, h) + T\beta(1 - \beta)^{T-1}[L(\mathcal{D}, h) + \ell_h^{\text{CS}}(h')] + [1 - (1 - \beta)^T - T\beta(1 - \beta)^{T-1}]c_{\max} \tag{7}$$

$$= L(\mathcal{D}, h) + T\beta(1 - \beta)^{T-1}\ell_h^{\text{CS}}(h') + [1 - (1 - \beta)^T - T\beta(1 - \beta)^{T-1}](c_{\max} - L(\mathcal{D}, h)) \tag{8}$$

$$\leq L(\mathcal{D}, h) + T\beta\ell_h^{\text{CS}}(h') + [1 - (1 - \beta)^T - T\beta(1 - \beta)^{T-1}]c_{\max} \tag{9}$$

$$= L(\mathcal{D}, h) + T\beta\ell_h^{\text{CS}}(h') + \left( \sum_{i=2}^T (-1)^i \beta^i \binom{T}{i} \right) c_{\max} \tag{10}$$

$$\leq L(\mathcal{D}, h) + T\beta\ell_h^{\text{CS}}(h') + \frac{1}{2}T^2\beta^2 c_{\max}. \tag{11}$$

The first inequality writes out the precise probability of the events in terms of  $\beta$  and bounds the loss of the last event ( $c > 2$ ) by  $c_{\max}$ . The second inequality is algebraic. The third uses the assumption that  $\beta < 1/T$ . □

This lemma states that applying a single iteration of SEARN does not cause the structured prediction loss of the learned hypothesis to degrade too much. In particular, up to a first order approximation, the loss increases proportional to the loss of the learned classifier. This observation can be iterated to yield the following lemma:

**Lemma 2 (Iteration)** *For all  $\mathcal{D}$ , for all learned  $h'$ , after  $C/\beta$  iterations of SEARN beginning with a policy  $\pi$  with loss  $L(\mathcal{D}, \pi)$ , and average learned losses as (4), the loss of the final learned policy  $h$  (without the optimal policy component) is bounded by (12):*

$$L(\mathcal{D}, h_{\text{last}}) \leq L(\mathcal{D}, \pi) + CT\ell_{\text{avg}} + c_{\max} \left( \frac{1}{2}CT^2\beta + T \exp[-C] \right). \tag{12}$$

This lemma states that after  $C/\beta$  iterations of SEARN the learned policy is not much worse than the quality of the initial policy  $\pi$ . The theorem follows from a choice of the constants  $\beta$  and  $C$  in Lemma 2.

*Proof* The proof involves invoking Lemma 1  $C/\beta$  times. The second and the third terms sum to give the following:

$$L(\mathcal{D}, h) \leq L(\mathcal{D}, \pi) + CT\ell_{\text{avg}} + c_{\max} \left( \frac{1}{2}CT^2\beta \right).$$

Last, if we call the initial policy, we fail with loss at most  $c_{\max}$ . The probability of failure after  $C/\beta$  iterations is at most  $T(1 - \beta)^{C/\beta} \leq T \exp[-C]$ . □

## 5 Comparison to alternative techniques

Standard techniques for structured prediction focus on the case where the  $\arg \max$  in (13) is tractable. Given its tractability, they attempt to learn parameters  $\theta$  such that solving (13) often results in low loss. There are a handful of classes of such algorithms and a large number of variants of each. Here, we focus on *independent classifier* models, *perceptron-based* models, and *global* models (such as conditional random fields and max-margin Markov networks). There are, of course, alternative frameworks (see, e.g., Weston et al. 2002; McAllester et al. 2004; Altun et al. 2004; McDonald et al. 2004; Tsochantaridis et al. 2005), but these are common examples.

### 5.1 The $\arg \max$ problem

Many structured prediction problems construct a scoring function  $F(y|x, \theta)$ . For a given input  $x \in \mathcal{X}$  and set of parameters  $\theta \in \Theta$ ,  $F$  provides a score for each possible output  $y$ . This leads to the “ $\arg \max$ ” problem (also known as the decoding problem or the pre-image problem), which seeks to find the  $y$  that maximizes  $F$  in order to make a prediction:

$$\hat{y} = \arg \max_{y \in \mathcal{Y}_x} F(y|x, \theta). \quad (13)$$

In (13), we seek the output  $y$  from the set  $\mathcal{Y}_x$  (where  $\mathcal{Y}_x \subseteq \mathcal{Y}$  is the set of all “reasonable” outputs for the input  $x$  – typically assumed to be finite). Unfortunately, solving (13) exactly is tractable only for very particular structures  $\mathcal{Y}$  and scoring functions  $F$ . As an easy example, when  $\mathcal{Y}_x$  is interpreted as a label sequence and the score function  $F$  depends only on adjacent labels, then dynamic programming can be used, leading to an  $\mathcal{O}(nk^2)$  prediction algorithm, where  $n$  is the length of the sequence and  $k$  is the number of possible labels for each element in the sequence. Similarly, if  $\mathcal{Y}$  represents trees and  $F$  obeys a context-free assumption, then this problem can be solved in time  $\mathcal{O}(n^3k)$ .

Often we are interested in more complex structures, more complex features or both. For such tasks, an exact solution to (13) is not tractable. For example, In natural language processing most statistical word-based and phrase-based models of translation are known to be NP-hard (Germann et al. 2003); syntactic translations models based on synchronous context free grammars are sometimes polynomial, but with an exponent that is too large in practice, such as  $n^{11}$  (Huang et al. 2005). Even in comparatively simple problems like sequence labeling and parsing—which are only  $\mathcal{O}(n)$  or  $\mathcal{O}(n^3)$ —it is often still computationally prohibitive to perform exhaustive search (Bikel 2004). For another sort of example, in computational biology, most models for phylogeny (Foulds and Graham 1982) and protein secondary structure prediction (Crescenzi et al. 1998) result in NP-hard search problems.

When faced with such intractable search problem, the standard tactic is to use an approximate search algorithm, such as greedy search, beam search, local hill-climbing search, simulated annealing, etc. These search algorithms are unlikely to be provably optimal (since this would imply that one is efficiently solving an NP-hard problem), but the hope is that they perform well on problems that are observed in the real world, as opposed to “worst case” inputs.

Unfortunately, applying suboptimal search algorithms to solve the structured prediction problem from (13) dispenses with many nice theoretical properties enjoyed by sophisticated learning algorithms. For instance, it may be possible to learn Bayes-optimal parameters  $\theta$  such that *if* exact search were possible, one would always find the best output. But given that exact search is not possible, such properties go away. Moreover, given that different

search algorithms exhibit different properties and biases, it is easy to believe that the value of  $\theta$  that is optimal for one search algorithm is not the same as the value that is optimal for another search algorithm.<sup>4</sup> It is these observations that have motivated our exploration of *search-based* structured prediction algorithms: learning algorithms for structured prediction that explicitly model the search process.

## 5.2 Independent classifiers

There are essentially two varieties of local classification techniques applied to structured prediction problems. In the first variety, the structure in the problem is ignored, and a single classifier is trained to predict each element in the output vector independently (Punyakanok and Roth 2001) or with dependence created by enforcement of membership in  $\mathcal{Y}_x$  constraints (Punyakanok et al. 2005b). The second variety is typified by maximum entropy Markov models (McCallum et al. 2000), though the basic idea of MEMMs has also been applied more generally to SVMs (Kudo and Matsumoto 2001, 2003; Giménez and Màrquez 2004). In this variety, the elements in the prediction vector are made sequentially, with the  $n$ th element conditional on outputs  $n - k \dots n - 1$  for a  $k$ th order model.

In the purely independent classifier setting, both training and testing proceed in the obvious way. Since the classifiers make one decision completely independently of any other decision, training makes use only of the input. This makes training the classifiers incredibly straightforward, and also makes prediction easy. In fact, running SEARN with  $\Phi(x, y)$  independent of all but  $y_n$  for the  $n$  prediction would yield exactly this framework (note that there would be no reason to iterate SEARN in this case). While this renders the independent classifiers approach attractive, it is also significantly weaker, in the sense that one cannot define complex features over the output space. This has not thus far hindered its applicability to problems like sequence labeling (Punyakanok and Roth 2001), parsing and semantic role labeling (Punyakanok et al. 2005a), but does seem to be an overly strict condition. This also limits the approach to Hamming loss.

SEARN is more similar to the MEMM-esque prediction setting. The key difference is that in the MEMM, the  $n$ th prediction is being made on the basis of the  $k$  previous predictions. However, these predictions are noisy, which potentially leads to the suboptimal performance described in the previous section. The essential problem is that the models have been trained assuming that they make all previous predictions correctly, but when applied in practice, they only have predictions about previous labels. It turns out that this can cause them to perform nearly arbitrarily badly. This is formalized in the following theorem, due to Matti Kääriäinen.

**Theorem 3** (Kääriäinen 2006) *There exists a distribution  $\mathcal{D}$  over first order binary Markov problems such that training a binary classifier based on true previous predictions to an error rate of  $\epsilon > 0$  leads to a Hamming loss given in (14), where  $T$  is the length of the sequence:*

$$\frac{T}{2} - \frac{1 - (1 - 2\epsilon)^{T+1}}{4\epsilon} + \frac{1}{2} \approx \frac{T}{2} \quad (14)$$

where the approximation is true for small  $\epsilon$  or large  $T$ .

<sup>4</sup>In fact, (Wainwright 2006) has provided evidence that when using approximation algorithms for graphical models, it is important to use the same approximate at both training and testing time.

Recently, Cohen and Carvalho (2005) has described an algorithm termed *stacked sequential learning* that attempts to remove this bias from MEMMs in a similar fashion to SEARN. The stacked algorithm learns a sequence of MEMMs, with the model trained on the  $(t + 1)$ st iteration based on outputs of the model from the  $t$ th iteration. For sequence labeling problems, this is quite similar to the behavior of SEARN when  $\beta$  is set to 1. However, unlike SEARN, the stacked sequential learning framework is effectively limited to sequence labeling problems. This limitation arises from the fact that it implicitly assumes that the set of decisions one must make in the future are always going to be same, regardless of decisions in the past. In many applications, such as entity detection and tracking (Daumé III and Marcu 2005b), this is not true. The set of possible choices (actions) available at time step  $i$  is heavily dependent on past choices. This makes the stacked sequential learning inapplicable in these problems.

### 5.3 Perceptron-style algorithms

The structured perceptron is an extension of the standard perceptron (Rosenblatt 1958) to structured prediction (Collins 2002). Assuming that the arg max problem is tractable, the structured perceptron constructs the weight vector in nearly an identical manner as for the binary case. While looping through the training data, whenever the predicted  $\hat{y}_n$  for  $x_n$  differs from  $y_n$ , we update the weights according to (15):

$$\mathbf{w} \leftarrow \mathbf{w} + \Phi(x_n, y_n) - \Phi(x_n, \hat{y}_n). \quad (15)$$

This weight update serves to bring the vector closer to the true output and further from the incorrect output. As in the standard perceptron, this often leads to a learned model that generalizes poorly. As before, one solution to this problem is weight averaging (Freund and Shapire 1999).

The *incremental perceptron* (Collins and Roark 2004) is a variant on the structured perceptron that deals with the issue that the arg max may not be analytically available. The idea of the incremental perceptron is to replace the arg max with a beam search algorithm. The key observation is that it is often possible to detect in the process of executing search whether it is possible for the resulting output to ever be correct. The incremental perceptron is essentially a search-based structured prediction technique, although it was initially motivated only as a method for speeding up convergence of the structured perceptron. In comparison to SEARN, it is, however, much more limited. It cannot cope with arbitrary loss functions, and is limited to a beam-search application. Moreover, for search problems with a large number of internal decisions (such as entity detection and tracking Daumé III and Marcu 2005b), aborting search at the first error is far from optimal.

### 5.4 Global prediction algorithms

Global prediction algorithms attempt to learn parameters that, essentially, rank correct (low loss) outputs higher than incorrect (high loss) alternatives.

Conditional random fields are an extension of logistic regression (maximum entropy models) to structured outputs (Lafferty et al. 2001). Similar to the structured perceptron, a conditional random field does not employ a loss function, but rather optimizes a log-loss approximation to the 0/1 loss over the entire output. Only when the features and structure are chosen properly can dynamic programming techniques be used to compute the required partition function, which typically limits the application of CRFs to linear chain models under a Markov assumption.

The maximum margin Markov network ( $M^3N$ ) formalism considers the structured prediction problem as a quadratic programming problem (Taskar et al. 2003, 2005), following the formalism for the support vector machine for binary classification. The  $M^3N$  formalism extends this to structured outputs under a given loss function  $l$  by requiring that the *difference* in score between the true output  $y$  and any incorrect output  $\hat{y}$  is at least the loss  $l(x, y, \hat{y})$  (modulo slack variables). That is: the  $M^3N$  framework scales the *margin* to be proportional to the loss. Under restrictions on the output space and the features (essentially, linear chain models with Markov features) it is possible to solve the corresponding quadratic program in polynomial time.

In comparison to CRFs and  $M^3Ns$ , SEARN is strictly more general. SEARN is limited neither to linear chains nor to Markov style features and can effectively and efficiently optimize structured prediction models under far weaker assumptions (see Sect. 6.2 for empirical evidence supporting this claim).

## 6 Experimental results

In this section, we present experimental results on two different sorts of structured prediction problems. The first set of problems—the sequence labeling problems—are comparatively simple and are included to demonstrate the application of SEARN to easy tasks. They are also the most common application domain on which other structured prediction techniques are tested; this enables us to directly compare SEARN with alternative algorithms on standardized data sets. The second application we describe is based on an automatic document summarization task, which is a significantly more complex domain than sequence labeling. This task enables us to test SEARN on significantly more complex problems with loss functions that do not decompose over the structure.

### 6.1 Sequence labeling

Sequence labeling is the task of assigning a *label* to each element in an input sequence. Sequence labeling is an attractive test bed for structured prediction algorithms because it is the simplest non-trivial structure. Modern state-of-the-art structured prediction techniques fare very well on sequence labeling problems. In this section, we present a range of results investigating the performance of SEARN on four separate sequence labeling tasks: handwriting recognition, named entity recognition (in Spanish), syntactic chunking and joint chunking and part-of-speech tagging.

For pure sequence labeling tasks (i.e., when segmentation is not also done), the standard loss function is Hamming loss, which gives credit on a per label basis. For a true output  $y$  of length  $N$  and hypothesized output  $\hat{y}$  (also of length  $N$ ), Hamming loss is defined according to (16):

$$\ell^{\text{Ham}}(y, \hat{y}) \triangleq \sum_{n=1}^N \mathbf{1}[y_n \neq \hat{y}_n]. \quad (16)$$

The most common loss function for joint segmentation and labeling problems (like the named entity recognition and syntactic chunking problems) is  $F_1$  measure over chunks.<sup>5</sup>  $F_1$

<sup>5</sup>We note in passing that directly optimizing  $F_1$  may not be the best approach, from the perspective of integrating information in a pipeline (Manning 2006). However, since  $F_1$  is commonly used and does *not* decompose over the output sequence, we use it for the purposes of demonstration.

is the geometric mean of precision and recall over the (properly-labeled) chunk identification task, given in (17):

$$e^F(y, \hat{y}) \triangleq \frac{2|y \cap \hat{y}|}{|y| + |\hat{y}|}. \quad (17)$$

As can be seen in (17), one is penalized both for identifying too many chunks (penalty in the denominator) and for identifying too few (penalty in the numerator). The advantage of  $F_1$  measure over Hamming loss seen most easily in problems where the majority of words are “not chunks”—for instance, in gene name identification (McDonald and Pereira 2005)—Hamming loss often prefers a system that identifies *no* chunks to one that identifies some correctly and other incorrectly. Using a weighted Hamming loss cannot completely alleviate this problem, for essentially the same reasons that a weighted zero-one loss cannot optimize  $F_1$  measure in binary classification, though one can often achieve an approximation (Lewis 2001; Musicant et al. 2003).

### 6.1.1 Handwriting recognition

The handwriting recognition task we consider was introduced by (Kassel 1995). Later, (Taskar et al. 2003) presented state-of-the-art results on this task using max-margin Markov networks. The task is an image recognition task: the input is a sequence of pre-segmented hand-drawn letters and the output is the character sequence (“a”-“z”) in these images. The data set we consider is identical to that considered by (Taskar et al. 2003) and includes 6600 sequences (words) collected from 150 subjects. The average word contains 8 characters. The images are  $8 \times 16$  pixels in size, and rasterized into a binary representation. Example image sequences are shown in Fig. 2 (the first characters are removed because they are capitalized).

For each possible output letter, there is a unique feature that counts how many times that letter appears in the output. Furthermore, for each pair of letters, there is an “edge” feature counting how many times this pair appears adjacent in the output. These edge features are the *only* “structural features” used for this task (i.e., features that span multiple output labels). Finally, for every output letter and for every pixel position, there is a feature that counts how many times that pixel position is “on” for the given output letter.

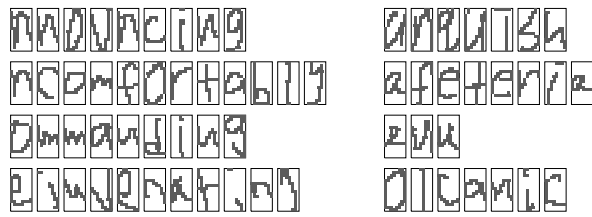
In the experiments, we consider two variants of the data set. The first, “small,” is the problem considered by (Taskar et al. 2003). In the small problem, ten fold cross-validation is performed over the data set; in each fold, roughly 600 words are used as training data and the remaining 6000 are used as test data. In addition to this setting, we also consider the “large” reverse experiment: in each fold, 6000 words are used as training data and 600 are used as test data.

### 6.1.2 Spanish named entity recognition

The named entity recognition (NER) task is concerned with spotting names of persons, places and organizations in text. Moreover, in NER we only aim to spot *names* and neither pronouns (“he”) nor nominal references (“the President”). We use the CoNLL 2002 data set, which consists of 8324 training sentences and 1517 test sentences; examples are shown in Fig. 3. A 300-sentence subset of the training data set was previously used by (Tsochantaridis et al. 2005) for evaluating the SVM<sup>struct</sup> framework in the context of sequence labeling. The small training set was likely used for computational considerations. The best reported results to date using the full data set are due to (Ando and Zhang 2005). We report results on both the “small” and “large” data sets.



**Fig. 2** Eight example words from the handwriting recognition data set



El presidente de la [Junta de Extremadura]<sub>ORG</sub> , [Juan Carlos Rodríguez Ibarra]<sub>PER</sub> , recibirá en la sede de la [Presidencia del Gobierno]<sub>ORG</sub> extremeño a familiares de varios de los condenados por el proceso “ [Lasa-Zabala]<sub>MISC</sub> ” , entre ellos a [Lourdes Díez Urraca]<sub>PER</sub> , esposa del ex gobernador civil de [Guipúzcoa]<sub>LOC</sub> [Julen Elgorriaga]<sub>PER</sub> ; y a [Antonio Rodríguez Galindo]<sub>PER</sub> , hermano del general [Enrique Rodríguez Galindo]<sub>PER</sub> .

**Fig. 3** Example labeled sentence from the Spanish Named Entity Recognition task

The structural features used for this task are roughly the same as in the handwriting recognition case. For each label, each label pair and each label triple, a feature counts the number of times this element is observed in the output. Furthermore, the standard set of input features includes the words and simple functions of the words (case markings, prefix and suffix up to three characters) within a window of  $\pm 2$  around the current position. These input features are paired with the current label. This feature set is fairly standard in the literature, though (Ando and Zhang 2005) report significantly improved results using a much larger set of features. In the results shown later in this section, all comparison algorithms use identical feature sets.

### 6.1.3 Syntactic chunking

The final sequence labeling task we consider is syntactic chunking (for English), based on the CoNLL 2000 data set. This data set includes 8936 sentences of training data and 2012 sentences of test data. An example is shown in Fig. 4. (Several authors have considered the *noun-phrase chunking* task instead of the full syntactic chunking task. It is important to notice the difference, though results on these two tasks are typically very similar, indicating that the majority of the difficulty is with noun phrases.)

We use the same set of features across all models, separated into “base features” and “meta features.” The base features apply to words individually, while meta features apply to entire chunks. The standard base features used are: the chunk length, the word (original, lower cased, stemmed, and original-stem), the case pattern of the word, the first and last 1, 2 and 3 characters, and the part of speech and its first character. We additionally consider membership features for lists of names, locations, abbreviations, stop words, etc. The meta features we use are, for any base feature  $b$ ,  $b$  at position  $i$  (for any sub-position of the chunk),  $b$  before/after the chunk, the entire  $b$ -sequence in the chunk, and any 2- or 3-gram tuple of  $b$ s in the chunk. We use a first order Markov assumption (chunk label only depends on the most recent previous label) and all features are placed on labels, not on transitions. In the results shown later in this section, some of the algorithms use a slightly different feature set. In particular, the CRF-based model uses similar, but not identical features; see (Sutton et al. 2005) for details.

[Great American]<sub>NP</sub> [said]<sub>VP</sub> [it]<sub>NP</sub> [increased]<sub>VP</sub> [its loan-loss reserves]<sub>NP</sub> [by]<sub>PP</sub> [\$ 93 million]<sub>NP</sub> [after]<sub>PP</sub> [reviewing]<sub>VP</sub> [its loan portfolio]<sub>NP</sub> , [raising]<sub>VP</sub> [its total loan and real estate reserves]<sub>NP</sub> [to]<sub>PP</sub> [\$ 217 million]<sub>NP</sub> .

Fig. 4 Example labeled sentence from the syntactic chunking task

Great<sup>NNP</sup><sub>B-NP</sub> American<sup>NNP</sup><sub>I-NP</sub> said<sup>VBD</sup><sub>B-VP</sub> it<sup>PRP</sup><sub>B-NP</sub> increased<sup>VBD</sup><sub>B-VP</sub> its<sup>PRP\$</sup><sub>B-NP</sub> loan-loss<sup>NN</sup><sub>I-NP</sub> reserves<sup>NNS</sup><sub>I-NP</sub> by<sup>IN</sup><sub>B-PP</sub> \$<sup>\$</sup><sub>B-NP</sub> 93<sup>CD</sup><sub>I-NP</sub> million<sup>CD</sup><sub>I-NP</sub> after<sup>IN</sup><sub>B-PP</sub> reviewing<sup>VBG</sup><sub>B-VP</sub> its<sup>PRP\$</sup><sub>B-NP</sub> loan<sup>NN</sup><sub>I-NP</sub> portfolio<sup>NN</sup><sub>I-NP</sub> .<sub>O</sub>

Fig. 5 Example sentence for the joint POS tagging and syntactic chunking task

### 6.1.4 Joint chunking and tagging

In the preceding sections, we considered the single sequence labeling task: to each element in a sequence, a single label is assigned. In this section, we consider the *joint* sequence labeling task. In this task, each element in a sequence is labeled with *multiple* tags. A canonical example of this task is joint POS tagging and syntactic chunking (Sutton et al. 2004). An example sentence jointly labeled for these two outputs is shown in Fig. 5 (under the BIO encoding).

For SEARN, there is little difference between standard sequence labeling and joint sequence labeling. We use the same data set as for the standard syntactic chunking task (Sect. 6.1.3) and essentially the same features. In order to model the fact that the two streams of labels are not independent, we decompose the problem into two parallel tagging tasks. First, the first POS label is determined, then the first chunk label, then the second POS label, then the second chunk label, etc. The only difference between the features we use in this task and the vanilla chunking task has to do the structural features. The structural features we use include the obvious Markov features on the individual sequences: counts of singleton, doubleton and tripleton POS and chunk tags. We also use “crossing sequence” features. In particular, we use counts of pairs of POS and chunk tags at the same time period as well as pairs of POS tags at time  $t$  and chunk tags at  $t - 1$  and vice versa.

### 6.1.5 Search and initial policies

The choice of “search” algorithm in SEARN essentially boils down to the choice of output vector representation, since, as defined, SEARN always operates in a left-to-right manner over the output vector. In this section, we describe vector representations for the output space and corresponding optimal policies for SEARN.

The most natural vector encoding of the sequence labeling problem is simply as itself. In this case, the search proceeds in a greedy left-to-right manner with one word being labeled per step. This search order admits some linguistic plausibility for many natural language problems. It is also attractive because (assuming unit-time classification) it scales as  $\mathcal{O}(NL)$ , where  $N$  is the length of the input and  $L$  is the number of labels, independent of the number of features or the loss function. However, this vector encoding is also highly biased, in the sense that it is perhaps not optimal for some (perhaps unnatural) problems. Other orders are possible (such as allowing any arbitrary position to be labeled at any time, effectively mimicing belief propagation); see (Daumé III 2006) for more experimental results under alternative orderings.

For joint segmentation and labeling tasks, such as named entity identification and syntactic chunking, there are two natural encodings: word-at-a-time and chunk-at-a-time. In

word-at-a-time, one essentially follows the “BIO encoding” and tags a single word in each search step. In chunk-at-a-time, one tags single *chunks* in each search step, which can consist of multiple words (after fixing a maximum phrase length). In our experiments, we focus exclusively on chunk-at-a-time decoding, as it is more expressive (feature-wise) and has been seen to perform better in other scenarios (Sarawagi and Cohen 2004).

Under the chunk-at-a-time encoding, an input of length  $N$  leads to a vector of length  $N$  over  $M \times L + 1$  labels, where  $M$  is the maximum phrase length. The interpretation of the first  $M \times L$  labels, for instance  $(m, l)$  means that the next phrase is of length  $m$  and is a phrase of type  $l$ . The “+1” label corresponds to a “complete” indicator. Any vector for which the sum of the “ $m$ ” components is not exactly  $N$  attains maximum loss.

### 6.1.6 Initial policies

For the sequence labeling problem under Hamming loss, the optimal policy is always to label the next word correctly. In the left-to-right order, this is straightforward. For the segmentation problem, word-at-a-time and chunk-at-a-time behave very similarly with respect to the loss function and optimal policy. We discuss word-at-a-time because its notationally more convenient, but the difference is negligible. The optimal policy can be computed by analyzing a few options in (18):

$$\pi(x, y_{1:T}, \hat{y}_{1:t-1}) = \begin{cases} \text{begin } X & y_t = \text{begin } X \\ \text{in } X & y_t = \text{in } X \text{ and } \hat{y}_{t-1} \in \{\text{begin } X, \text{in } X\}, \\ \text{out} & \text{otherwise.} \end{cases} \quad (18)$$

It is easy to show that this policy is optimal (assuming noise-free training data). There is, however, another equally optimal policy. For instance, if  $y_t$  is “in  $X$ ” but  $\hat{y}_{t-1}$  is “in  $Y$ ” (for  $X \neq Y$ ), then it is equally optimal to select  $\hat{y}_t$  to be “out” or “in  $Y$ ”. In theory, when the optimal policy does not care about a particular decision, one can randomize over the selection. However, in practice, we always default to a particular choice to reduce noise in the learning process.

For all of the policies described above, it is also straightforward to compute the optimal approximation for estimating the expected cost of an action. In the Hamming loss case, the loss is 0 if the choice is correct and 1 otherwise. The computation for  $F_1$  loss is a bit more complicated: one needs to compute an optimal intersection size for the future and add it to the past “actual” size. This is also straightforward by analyzing the same cases as in (18).

### 6.1.7 Experimental results and discussion

In this section, we compare the performance of SEARN to the performance of alternative structured prediction techniques over the data sets described above. The results of this evaluation are shown in Table 1. In this table, we compare raw classification algorithms (perceptron, logistic regression and SVMs) to alternative structured prediction algorithms (structured perceptron, CRFs, SVM<sup>struct</sup>s and M<sup>3</sup>Ns) to SEARN with three baseline classifiers (perceptron, logistic regression and SVMs). For all SVM algorithms and for M<sup>3</sup>Ns, we compare both linear and quadratic kernels (cubic kernels were evaluated but did not lead to improved performance over quadratic kernels).

For all SEARN-based models, we use the following settings of the tunable parameters (see Daumé III 2006 for a comparison of different settings). We use the optimal approximation for the computation of the per-action costs. We use a left-to-right search order with a beam of size 10. For the chunking tasks, we use chunk-at-a-time search. We use weighted all pairs

**Table 1** Empirical comparison of performance of alternative structured prediction algorithms against SEARN on sequence labeling tasks. (Top) Comparison for whole-sequence 0/1 loss; (Bottom) Comparison for individual losses: Hamming for handwriting and Chunking+Tagging and F for NER and Chunking. SEARN is always optimized for the appropriate loss

ALGORITHM	Handwriting		NER		Chunk	C+T
	Small	Large	Small	Large		
<b>CLASSIFICATION</b>						
Perceptron	65.56	70.05	91.11	94.37	83.12	87.88
Log Reg	68.65	72.10	93.62	96.09	85.40	90.39
SVM-Lin	75.75	82.42	93.74	97.31	86.09	93.94
SVM-Quad	82.63	82.52	85.49	85.49	~	~
<b>STRUCTURED</b>						
Str. Perc.	69.74	74.12	93.18	95.32	92.44	93.12
CRF	–	–	94.94	~	94.77	96.48
SVM <sup>struct</sup>	–	–	94.90	~	–	–
M <sup>3</sup> N-Lin	81.00	~	–	–	–	–
M <sup>3</sup> N-Quad	87.00	~	–	–	–	–
<b>SEARN</b>						
Perceptron	70.17	76.88	95.01	97.67	94.36	96.81
Log Reg	73.81	79.28	95.90	98.17	94.47	96.95
SVM-Lin	82.12	90.58	95.91	98.11	94.44	96.98
SVM-Quad	87.55	90.91	89.31	90.01	~	~

(Beygelzimer et al. 2005) and costing (Zadrozny et al. 2003) to reduce from cost-sensitive classification to binary classification.

Note that some entries in Table 1 are missing. The vast majority of these entries are missing because the algorithm considered could not reasonably scale to the data set under consideration. These are indicated with a “~” symbol. Other entries are not available simply because the results we report are copied from other publications and these publications did not report all relevant scores. These are indicated with a “–” symbol.

It should be noted that there are several issues that can account for a “~” symbol. In the case of the quadratic SVMs applied to **Chunk** and **C+T**, the issue is primarily that the number of examples is simply too large to solve with libSVM (Chang and Lin 2001), which we used in our implementation. It is possible (in fact, likely) that some newer subgradient-based optimizer may be able to scale sufficiently. The same holds for the M<sup>3</sup>N results on the large **Handwriting** task. However, we do not believe this comparison to be unfair: the SEARN-based algorithms could equally benefit from such improved optimization methods. The second issue is that for **NER** under CRFs and SVM<sup>struct</sup>s, the inference is over segmentations, not labellings. This is considerably more expensive than simple labeling tasks. However, there is an additional issue: the time complexity for a gradient computation for a CRF is *exponential* in the Markov length; while it is independent of the Markov length for SEARN. (In addition to the big-O complexity, running the forward-backward algorithm is simply an expensive computation: the constant is large.)

We observe several patterns in the results from Table 1. The first is that structured techniques consistently outperform their classification counterparts (e.g., CRFs outperform logistic regression). The single exception is on the small handwriting task: the quadratic SVM

outperforms the quadratic  $M^3N$ .<sup>6</sup> For all classifiers, adding SEARN consistently improves performance.

An obvious pattern worth noticing is that moving from the small data set to the large data set results in improved performance, regardless of learning algorithm. However, equally interesting is that simple classification techniques when applied to large data sets outperform complicated learning techniques applied to small data sets. Although this comparison is not completely fair—both algorithms should get access to the same data—if the algorithm (like the SVM<sup>struct</sup> or the  $M^3N$ ) cannot scale to the large data set, then something is missing. For instance, a vanilla SVM on the large handwriting data set outperforms the  $M^3N$  on the small set. Similarly, a vanilla logistic regression classifier trained on the large NER data set outperforms the SVM<sup>struct</sup> and the CRF on the small data sets.

On the same data set, SEARN can perform comparable or better than competing structured prediction techniques. On the small handwriting task, the two best performing systems are  $M^3N$ s with quadratic kernels (87.0% accuracy) and SEARN with quadratic SVMs (87.6% accuracy). On the NER task, SEARN with a perceptron classifier performs comparably to SVM<sup>struct</sup> and CRFs (at around 95.9% accuracy). On the Chunking+Tagging task, all varieties of SEARN perform comparatively to the CRF. In fact, the only task on which SEARN does *not* outperform the competing techniques is on the raw chunking task, for which the CRF obtains an F-score of 94.77 compared to 94.47 for SEARN, using a significantly different feature set.

The final result from Table 1 worth noticing is that, with the exception of the handwriting recognition task, SEARN using logistic regression as a base learner performs at the top of the pack. The SVM-based SEARN models typically perform slightly better, but not significantly. In fact, the raw averaged perceptron with SEARN performs almost as well as the logistic regression. This is a nice result because the SVM-based models tend to be expensive to train, especially in comparison to the perceptron. The fact that this pattern does not hold for the handwriting task is likely due to the fact that the data for this task is quite unlike the data for the other tasks. For the handwriting task, there are a comparatively small number of features which are individually much less predictive of the class. It is only in combination that good classifiers can be learned.

While these results are useful, they should be taken with a grain of salt. Sequence labeling is a very easy problem. The structure is simple and the most common loss functions decompose over the structure. The comparatively good performance of raw classifiers suggests that the importance of structure is minor. In fact, some results suggest that one need not actually consider the structure at all for some such problems (Punyakanok and Roth 2001; Punyakanok et al. 2005b).

### 6.1.8 Effect of search order

One significant issue with SEARN, or any search-based prediction algorithm, is the choice of the prediction order. In the sequence labeling experiments presented in this section, the search order has always been left-to-right. This is sensible, as it is how humans process text. However, the question arises: could one do better? For certain problems, the answer seems to be “yes.” For certain sequence labeling tasks, it has been observed that ordering the search based on the *entropy* of the label distribution of certain words can lead to improved performance (Tsuruoka and Tsujii 2005). It was later observed that the order could be *learned*

<sup>6</sup>However, it should be noted that a different implementation technique was used in this comparison. The  $M^3N$  is based on an SMO algorithm, while the quadratic SVM is libsvm (Chang and Lin 2001).

(Shen et al. 2007). In the context of incremental parsing, for some languages, it is helpful to parse *backwards* rather than forwards (Turian and Melamed 2006). To ascertain the effect of search order on the NER task, we have run additional experiments in a right-to-left setting. In comparison to the left-to-right performance of 97.67 (for perceptron) and 98.17 (for logistic regression), the right-to-left performance was slightly worse: 97.40 and 97.98, respectively. These differences were not statistically significant at the 5% level.

A more substantial question is: what happens when there is *not* an obvious search order. This arises, for instance, in image segmentation, or collective-classification problems. One answer, akin to the learned search order method (Shen et al. 2007) would be to simply allow SEARN to choose the order on its own. This makes inference significantly more computationally intensive (but still much less so than exact probabilistic inference on the associated graph). This is the approach we take in the automatic document summarization example described in Sect. 6.2 where a left-to-right search order is not natural.

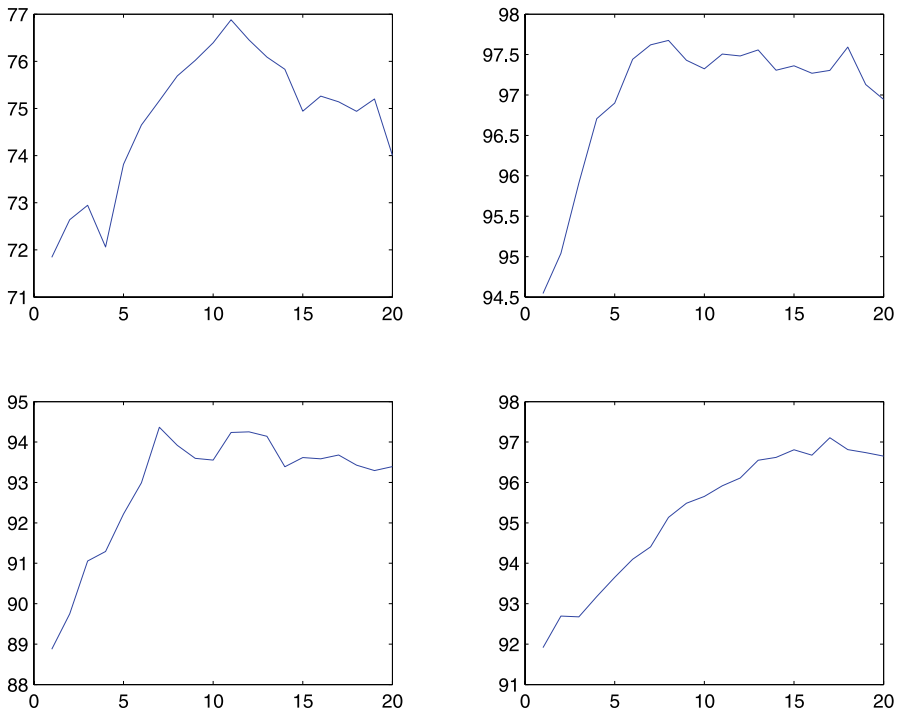
### 6.1.9 Effect of iterating

In Fig. 6, we plot the performance of the learned policy for the SEARN-based models for the four tasks as the number of iterations increases. For these graphs, we use a constant value of  $\beta = 1$  for the interpolation: pure policy-iteration. The curves are somewhat different for each problem, but in general an optimum is reached in 5–15 iterations and then performance either levels off (e.g., for syntactic chunking) or begins to drop (e.g., for handwriting recognition). The drop in performance is likely due to overfitting. Note that these curves are the performance of the learned policy *without* the optimal policy on the test data, so these graphs do not contradict the SEARN theorem of uniform degradation of performance (Lemma 1).

## 6.2 Automatic document summarization

Multidocument summarization is the task of creating a summary out of a collection of documents on a focused topic. In query-focused summarization, this topic is given explicitly in the form of a user’s query. The dominant approach to the multidocument summarization problem is sentence extraction: a summary is created by greedily extracting sentences from the document collection until a pre-defined word limit is reached. Teufel and Moens (1997) and Lin and Hovy (2002) describe representative examples. Recent work in sentence compression (Knight and Marcu 2002; McDonald 2006) and document compression (Daumé III and Marcu 2002) attempts to take small steps beyond sentence extraction. Compression models can be seen as techniques for extracting sentences then dropping extraneous information. They are more powerful than simple sentence extraction systems, while remaining trainable and tractable. Unfortunately, their training hinges on the existence of (sentence, compression) pairs, where the compression is obtained from the sentence by only dropping words and phrases (the work of Turner and Charniak 2005 is an exception). Obtaining such data is quite challenging.

The exact model we use for the document summarization task is a novel “vine-growth” model, described in more detail in (Daumé III 2006). The vine-growth method uses syntactic parses of the sentence in the form of *dependency* structures. In the vine-growth model, if a word  $w$  is to be included in the summary, then all words closer to the tree root are included. Viewed as a graph (for instance, in a CRF application), this would result in a (nearly) fully connected graph, for which we know inference is very hard.



**Fig. 6** Number of iterations of SEARN for each of the four sequence labeling problem. *Upper-left*: Handwriting recognition; *Upper-right*: Spanish named entity recognition; *Lower-left*: Syntactic chunking; *Lower-right*: Joint chunking/tagging

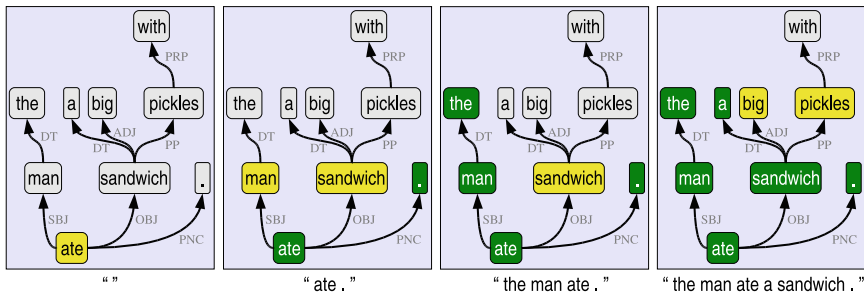
### 6.2.1 Search space and actions

The search algorithm we employ for implementing the vine-growth model is based on incrementally *growing* summaries. In essence, beginning with an empty summary, the algorithm incrementally adds words to the summary, either by beginning a new sentence or growing existing sentences. At any step in search, the root of a new sentence may be added, as may any direct child of a previously added node. To see more clearly how the vine-growth model functions, consider Fig. 7. This figure shows a four step process for creating the summary “the man ate a sandwich.” from the original document sentence “the man ate a big sandwich with pickles.”

When there is more than one sentence in the source documents, the search proceeds asynchronously across all sentences. When the sentences are laid out adjacently, the end summary is obtained by taking all the green summary nodes once a pre-defined word limit has been reached. This final summary is a collection of subtrees grown off a sequence of underlying trees: hence the name “vine-growth.”

### 6.2.2 Data and evaluation criteria

For data, we use the DUC 2005 data set (Dang 2005). This consists of 50 document collections of 25 documents each; each document collection includes a human-written query. Each document collection additionally has five human-written “reference” summaries (250



**Fig. 7** An example of the creation of a summary under the vine-growth model

words long, each) that serve as the gold standard. In the official DUC evaluations, all 50 collections are “test data.” However, since the DUC 2005 task is significantly different from previous DUC tasks, there is no a good source of training data. Therefore, we report results based on 10-fold cross validation. We train on 45 collections and test on the remaining 5.

Automatic evaluation is a notoriously difficult problem for document summarization. The current popular choice for metric is Rouge (Lin and Hovy 2003), which (roughly speaking) computes *n*-gram overlap between a system summary and a set of human written summaries. In various experiments, Rouge has been seen to correspond with human judgment of summary quality. In the experiments described in this chapter, we use the “Rouge 2” metric, which uses evenly weighted bigram scores.

### 6.2.3 Initial policy

Computing the best label completion under Rouge metric for the vine-growth model is intractable. The intractability stems from the model constraint that a word can only be added to a summary after its parent is added. We therefore use an approximate, search-based policy (see Sect. 3.4.2). In order to approximate the cost of a given partial summary, we *search* for the best possible completion. That is, if our goal is a 100 word summary and we have already created a 50 word summary, then we execute beam search (beam size 20) for the remaining 50 words that maximize the Rouge score.

### 6.2.4 Feature functions

Features in the vine-growth model may consider any aspect of the currently generated summary, and any part of the input document set. These features include simple lexical features: word identity, stem and part of speech of the word under consideration, the syntactic relation with its parent, the position and length of the sentence it appears in, whether it appears in quotes, the length of the document it appears in, the number of pronouns and attribution verbs in the subtree rooted at the word. The features also include language model probabilities for: the word, sentence and subtree under language models derived from the query, a BAYESUM representation of the query, and the existing partial summary.

### 6.2.5 Experimental results

Experimental results are shown in Table 2. We report Rouge scores for summaries of length 100 and length 250. We compare the following systems. First, oracle systems that perform the summarization task *with* knowledge of the true output, attempting to maximize



**Table 2** Summarization results; values are Rouge 2 scores (higher is better)

	ORACLE		SEARN		BAYESUM		Base	Best
	Vine	Extr	Vine	Extr	D05	D03		
100 w	0.0729	0.0362	0.0415	0.0345	0.0340	0.0316	0.0181	–
250 w	0.1351	0.0809	0.0824	0.0767	0.0762	0.0698	0.0403	0.0725

the Rouge score. We present results for an oracle sentence extraction system (Extr) and an oracle vine-growth system (Vine). Second, we present the results of the SEARN-based systems, again for both sentence extraction (Extr) and vine-growth (Vine). Both of these are trained with respect to the oracle system. (Note that it is impossible to compare against competing structured prediction techniques. This summarization problem, even in its simplified form, is far too complex to be amenable to other methods.) For comparison, we present results from the BAYESUM system (Daumé III and Marcu 2005a, 2006), which achieved the highest score according to human evaluations of responsiveness in DUC 05. This system, as submitted to DUC 05, was *trained* on DUC 2003 data; the results for this configuration are shown in the “D03” column. For the sake of fair comparison, we also present the results of this system, trained in the same cross-validation approach as the SEARN-based systems (column “D05”). Finally, we present the results for the baseline system and for the best DUC 2005 system (according to the Rouge 2 metric).

As we can see from Table 2 at the 100 word level, sentence extraction is a nearly solved problem for this domain and this evaluation metric. That is, the oracle sentence extraction system yields a Rouge score of 0.0362, compared to the score achieved by the SEARN system of 0.0345. This difference is on the border of statistical significance at the 95% level. The next noticeable item in the results is that, although the SEARN-based *extraction* system comes quite close to the theoretical optimal, the oracle results for the *vine-growth* method are significantly higher. Not surprisingly, under SEARN, the summaries produced by the vine-growth technique are uniformly better than those produced by raw extraction. The last aspect of the results to notice is how the SEARN-based models compare to the best DUC 2005 system, which achieved a Rouge score of 0.0725. The SEARN-based systems uniformly dominate this result, but this comparison is not fair due to the training data. We can approximate the expected improvement for having the new training data by comparing the BAYESUM system when trained on the DUC 2005 and DUC 2003 data: the improvement is 0.0064 absolute. When this result is added to the best DUC 2005 system, its score rises to 0.0789, which is better than the SEARN-based extraction system but not as good as the vine-growth system. It should be noted that the best DUC 2005 system *was* a purely extractive system (Ye et al. 2005).

## 7 Discussion and conclusions

In this paper, we have:

- Presented an algorithm, SEARN, for solving complex structured prediction problems with minimal assumptions on the structure of the output and loss function.
- Compared the performance of SEARN against standard structured prediction algorithms on standard sequence labeling tasks, showing that it is competitive with existing techniques.

- Described a novel approach to summarization—the vine-growth method—and applied SEARN to the underlying learning problem, yielding state-of-the-art performance on standardized summarization data sets.

There are many lenses through which one can view the SEARN algorithm.

From an applied perspective, SEARN is an easy technique for training models for which complex search algorithms must be used. For instance, when using multiclass logistic regression as a base classifier for Hamming loss, the first iteration of SEARN is identical to training a maximum entropy Markov model. The subsequent iterations of SEARN can be seen as attempting to get around the fact that MEMMs are trained *assuming* all previous decisions are made correctly. This assumption is false, of course, in practice. Similar recent algorithms such a decision-tree-based parsing (Turian and Melamed 2006) and perceptron-based machine translation (Liang et al. 2006) can also be seen as running a (slightly modified) first iteration of SEARN.

SEARN contrasts with more typical algorithms such as CRFs and M<sup>3</sup>Ns based on considering how information is shared at test time. Standard algorithms use exact (typically Viterbi) search to share full information across the entire output, “trading off” one decision for another. SEARN takes an alternative approach: it attempts to share information at *training time*. In particular, by training the classifier using a loss based on both past experience and future expectations, the training attempts to integrate this information during learning. This is not unsimilar to the “alternative objective” proposed by (Kakade et al. 2002) for CRFs. One approach is not necessarily better than the other; they are simply different ways to accomplish the same goal.

One potential limitation to SEARN is that when one trains a new classifier on the output of a previous iteration’s classifier, it is usually going to be the case that previous iteration’s classifier performs better on the training data than on the test data. This means that, although training via SEARN is likely preferable to training against *only* an initial policy, it can still be overly optimistic. Based on the experimental evidence, it appears that this has yet to be a serious concern, but it remains worrisome. There are two easy ways to combat this problem. The first is simply to attempt to ensure that the learned classifiers do not overfit at all. In practice, however, this can be difficult. Another approach with a high computational cost is cross-validation. Instead of training one classifier in each SEARN step, one could train ten, each holding out a different 10% of the data. When asked to run the “current” classifier on an example, the classifier not trained on the example is used. This does not completely remove the possibility of overfitting, but significantly lessens its likelihood.

A second limitation, pointed out by (Zhang 2006), is that there is a slight disparity between what SEARN does at a theoretical level and how SEARN functions in practice. In particular, SEARN does not *actually* start with the optimal policy. Even when we can compute the initial policy exactly, the “true outputs” on which this initial policy are based are potentially noisy. This means that while  $\pi$  is optimal for the *noisy* data, it is not optimal for the *true* data distribution. In fact, it is possible to construct noisy distributions where SEARN performs poorly.<sup>7</sup> Finding other initial policies which are closer to optimal in these situations is an open problem.

SEARN obeys a desirable theoretical property: given a good classification algorithm, one is guaranteed a good structured prediction algorithm. Importantly, this result is independent

<sup>7</sup>One can construct such a noisy distribution as follows. Suppose there is fundamental noise and a “safe” option which results in small loss. Suppose this safe option is always more than a one step deviation from the highly noisy “optimal” sequence. SEARN can be confused by this divergence.

of the size of the search space or the tractability of the search method. This shows that local learning—when done properly—can lead to good global performance. From the perspective of applied machine learning, SEARN serves as an interpreter through which engineers can easily make use of state-of-the-art machine learning techniques.

In the context of structured prediction algorithms, SEARN lies somewhere between global learning algorithms, such as  $M^3Ns$  and CRFs, and local learning algorithms, such as those described (Punyakanok and Roth 2001). The key difference between SEARN and global algorithms is in how uncertainty is handled. In global algorithms, the search algorithm is used at test time to propagate uncertainty across the structure. In SEARN, the prediction costs are used during training time to propagate uncertainty across the structure. Both contrast with local learning, in which no uncertainty is propagated.

From a wider machine learning perspective, SEARN makes more apparent the connection between reinforcement learning and structured prediction. In particular, structured prediction can be viewed as a reinforcement learning problem in a degenerate world in which all observations are available at the initial time step. However, there are clearly alternative middle-grounds between pure structured prediction and full-blown reinforcement learning (and natural applications—such as planning—in this realm) for which this connection might serve to be useful.

Despite these successes, there is much future work that is possible. One significant open question on the theoretical side is that of sample complexity: “How many examples do we need in order to achieve learning under additional assumptions?” Related problems of semi-supervised and active learning in the SEARN framework are also interesting and likely to produce powerful extensions. Another vein of research is in applying SEARN to domains other than language. Structured prediction problems arise in a large variety of settings (vision, biology, system design, compilers, etc.). For each of these domains, different sorts of search algorithms and different sorts of features are necessary. Although SEARN has been discussed largely as a method for solving structured prediction problems, it is, more generally, a method for integrating *search* and *learning*. This leads to potential applications of SEARN that fall strictly outside the scope of structured prediction.

## References

- Altun, Y., Hofmann, T., & Smola, A. (2004). Gaussian process classification for segmenting and annotating sequences. In *Proceedings of the international conference on machine learning (ICML)*.
- Ando, R., & Zhang, T. (2005). A framework for learning predictive structures from multiple tasks and unlabeled data. *Journal of Machine Learning Research*, 6, 1817–1853.
- Bagnell, J. A., Kakade, S., Ng, A., & Schneider, J. (2003). Policy search by dynamic programming. In *Neural information processing systems* (Vol. 16). Cambridge: MIT Press.
- Beygelzimer, A., Dani, V., Hayes, T., Langford, J., & Zadrozny, B. (2005). Error limiting reductions between classification tasks. In *Proceedings of the international conference on machine learning (ICML)*.
- Bikel, D. M. (2004). Intricacies of Collins’ parsing model. *Computational Linguistics*, 30(4), 479–511.
- Chang, C.-C., & Lin, C.-J. (2001). *LIBSVM: a library for support vector machines*. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- Cohen, W. W., & Carvalho, V. (2005). Stacked sequential learning. In *Proceedings of the international joint conference on artificial intelligence (IJCAI)*.
- Collins, M. (2002). Discriminative training methods for hidden Markov models: theory and experiments with perceptron algorithms. In *Proceedings of the conference on empirical methods in natural language processing (EMNLP)*.
- Collins, M., & Roark, B. (2004). Incremental parsing with the perceptron algorithm. In *Proceedings of the conference of the association for computational linguistics (ACL)*.
- Crescenzi, P., Goldman, D., Papadimitriou, C., Piccolboni, A., & Yannakakis, M. (1998). On the complexity of protein folding. In *ACM symposium on theory of computing (STOC)* (pp. 597–603).

- Dang, H. (Ed.). (2005). *Fifth document understanding conference (DUC-2005)*, Ann Arbor, MI, June 2005.
- Daumé III, H. (2006). *Practical structured learning for natural language processing*. PhD thesis, University of Southern California.
- Daumé III, H., & Marcu, D. (2002). A noisy-channel model for document compression. In *Proceedings of the conference of the association for computational linguistics (ACL)* (pp. 449–456).
- Daumé III, H., & Marcu, D. (2005a). Bayesian summarization at DUC and a suggestion for extrinsic evaluation. In *Document understanding conference*.
- Daumé III, H., & Marcu, D. (2005b). A large-scale exploration of effective global features for a joint entity detection and tracking model. In *Proceedings of the joint conference on human language technology conference and empirical methods in natural language processing (HLT/EMNLP)* (pp. 97–104).
- Daumé III, H., & Marcu, D. (2006). Bayesian query-focused summarization. In *Proceedings of the conference of the association for computational linguistics (ACL)*, Sydney, Australia.
- Foulds, L. R., & Graham, R. L. (1982). The Steiner problem in phylogeny is NP-complete. *Advances in Applied Mathematics*, 3, 43–49.
- Freund, Y., & Shapire, R. E. (1999). Large margin classification using the perceptron algorithm. *Machine Learning*, 37(3), 277–296.
- Germann, U., Jahr, M., Knight, K., Marcu, D., & Yamada, K. (2003). Fast decoding and optimal decoding for machine translation. *Artificial Intelligence*, 154(1–2), 127–143.
- Giménez, J., & Màrquez, L. (2004). SVMTool: a general POS tagger generator based on support vector machines. In *Proceedings of the 4th LREC*.
- Huang, L., Zhang, H., & Gildea, D. (2005). Machine translation as lexicalized parsing with hooks. In *Proceedings of the 9th international workshop on parsing technologies (IWPT-05)*, October 2005.
- Kääriäinen, M. (2006). Lower bounds for reductions. In *The atomic learning workshop (TTI-C)*, March 2006.
- Kakade, S., & Langford, J. (2002). Approximately optimal approximate reinforcement learning. In *Proceedings of the international conference on machine learning (ICML)*.
- Kakade, S., Teh, Y. W., & Roweis, S. (2002). An alternate objective function for Markovian fields. In *Proceedings of the international conference on machine learning (ICML)*.
- Kassel, R. (1995). *A comparison of approaches to on-line handwritten character recognition*. PhD thesis, Massachusetts Institute of Technology, Spoken Language Systems Group.
- Knight, K., & Marcu, D. (2002). Summarization beyond sentence extraction: a probabilistic approach to sentence compression. *Artificial Intelligence*, 139(1).
- Kudo, T., & Matsumoto, Y. (2001). Chunking with support vector machines. In *Proceedings of the conference of the North American chapter of the association for computational linguistics (NAACL)*.
- Kudo, T., & Matsumoto, Y. (2003). Fast methods for kernel-based text analysis. In *Proceedings of the conference of the association for computational linguistics (ACL)*.
- Lafferty, J., McCallum, A., & Pereira, F. (2001). Conditional random fields: probabilistic models for segmenting and labeling sequence data. In *Proceedings of the international conference on machine learning (ICML)*.
- Langford, J., & Zadrozny, B. (2005). Relating reinforcement learning performance to classification performance. In *Proceedings of the international conference on machine learning (ICML)*.
- Lewis, D. (2001). Applying support vector machines to the TREC-2001 batch filtering and routing tasks. In *Proceedings of the conference on research and developments in information retrieval (SIGIR)*.
- Liang, P., Bouchard-Côté, A., Klein, D., & Taskar, B. (2006). An end-to-end discriminative approach to machine translation. In *Proceedings of the joint international conference on computational linguistics and association of computational linguistics (COLING/ACL)*.
- Lin, C.-Y., & Hovy, E. (2002). From single to multi-document summarization: a prototype system and its evaluation. In *Proceedings of the conference of the association for computational linguistics (ACL)*, July 2002.
- Lin, C.-Y., & Hovy, E. (2003). Automatic evaluation of summaries using  $n$ -gram co-occurrence statistics. In *Proceedings of the conference of the North American chapter of the association for computational linguistics and human language technology (NAACL/HLT)*, Edmonton, Canada, 27 May–1 June 2003.
- Manning, C. (2006). *Doing named entity recognition? Don't optimize for F<sub>1</sub>*. Post on the NLPers Blog, 25 August 2006. <http://nlpers.blogspot.com/2006/08/doing-named-entity-recognition-dont.html>.
- McAllester, D., Collins, M., & Pereira, F. (2004). Case-factor diagrams for structured probabilistic modeling. In *Proceedings of the conference on uncertainty in artificial intelligence (UAI)*.
- McCallum, A., Freitag, D., & Pereira, F. (2000). Maximum entropy Markov models for information extraction and segmentation. In *Proceedings of the international conference on machine learning (ICML)*.
- McDonald, R. (2006). Discriminative sentence compression with soft syntactic constraints. In *Proceedings of the conference of the European association for computational linguistics (EACL)*.
- McDonald, R., & Pereira, F. (2005). Identifying gene and protein mentions in text using conditional random fields. *BMC Bioinformatics*, 6(Suppl 1).

- McDonald, R., Crammer, K., & Pereira, F. (2004). Large margin online learning algorithms for scalable structured classification. In *NIPS workshop on learning with structured outputs*.
- Musican, D., Kumar, V., & Ozgur, A. (2003). Optimizing F-measure with support vector machines. In *Proceedings of the international Florida artificial intelligence research society conference* (pp. 356–360).
- Ng, A., & Jordan, M. (2000). PEGASUS: A policy search method for large MDPs and POMDPs. In *Proceedings of the conference on uncertainty in artificial intelligence (UAI)*.
- Punyakanok, V., & Roth, D. (2001). The use of classifiers in sequential inference. In *Advances in neural information processing systems (NIPS)*.
- Punyakanok, V., Roth, D., & Yih, W.-T. (2005a). The necessity of syntactic parsing for semantic role labeling. In *Proceedings of the international joint conference on artificial intelligence (IJCAI)* (pp. 1117–1123).
- Punyakanok, V., Roth, D., Yih, W.-T., & Zimak, D. (2005b). Learning and inference over constrained output. In *Proceedings of the international joint conference on artificial intelligence (IJCAI)* (pp. 1124–1129).
- Rosenblatt, F. (1958). The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological Review*, 65, 386–408. Reprinted in *Neurocomputing* (MIT Press, 1998).
- Russell, S., & Norvig, P. (1995). *Artificial intelligence: a modern approach*. New Jersey: Prentice Hall.
- Sarawagi, S., & Cohen, W. (2004). Semi-Markov conditional random fields for information extraction. In *Advances in neural information processing systems (NIPS)*.
- Shen, L., Satta, G., & Joshi, A. (2007). Guided learning for bidirectional sequence classification. In *Proceedings of the conference of the association for computational linguistics (ACL)*.
- Sutton, C., Rohanimanesh, K., & McCallum, A. (2004). Dynamic conditional random fields: factorized probabilistic models for labeling and segmenting sequence data. In *Proceedings of the international conference on machine learning (ICML)* (pp. 783–790).
- Sutton, C., Sindelar, M., & McCallum, A. (2005). *Feature bagging: preventing weight undertraining in structured discriminative learning* (Technical Report IR-402). University of Massachusetts, Center for Intelligent Information Retrieval.
- Taskar, B., Guestrin, C., & Koller, D. (2003). Max-margin Markov networks. In *Advances in neural information processing systems (NIPS)*.
- Taskar, B., Chatalbashev, V., Koller, D., & Guestrin, C. (2005). Learning structured prediction models: a large margin approach. In *Proceedings of the international conference on machine learning (ICML)* (pp. 897–904).
- Teufel, S., & Moens, M. (1997). Sentence extraction as a classification task. In *ACL/EACL-97 workshop on intelligent and scalable text summarization* (pp. 58–65).
- Tsochantaridis, I., Hofmann, T., Joachims, T., & Altun, Y. (2005). Large margin methods for structured and interdependent output variables. *Journal of Machine Learning Research*, 6, 1453–1484.
- Tsuruoka, Y., & Tsujii, J. (2005). Bidirectional inference with the easiest-first strategy for tagging sequence data. In *Proceedings of the conference on empirical methods in natural language processing (EMNLP)*.
- Turian, J., & Melamed, I. D. (2006). Advances in discriminative parsing. In *Proceedings of the joint international conference on computational linguistics and association of computational linguistics (COLING/ACL)*.
- Turner, J., & Charniak, E. (2005). Supervised and unsupervised learning for sentence compression. In *Proceedings of the conference of the association for computational linguistics (ACL)*.
- Wainwright, M. (2006). *Estimating the “wrong” graphical model: benefits in the computation-limited setting* (Technical report). University of California Berkeley, Department of Statistics, February 2006.
- Weston, J., Chapelle, O., Elisseeff, A., Schoelkopf, B., & Vapnik, V. (2002). Kernel dependency estimation. In *Advances in neural information processing systems (NIPS)*.
- Ye, S., Qiu, L., Chua, T.-S., & Kan, M.-Y. (2005). NUS at DUC 2005: understanding documents via concept links. In *Document understanding conference*.
- Zadrozny, B., Langford, J., & Abe, N. (2003). Cost-sensitive learning by cost-proportionate example weighting. In *Proceedings of the IEEE conference on data mining (ICDM)*.
- Zhang, T. (2006). Personal communication, June 2006.