

4

Entropy-Guided Structure Learning Framework

We show in this chapter that the proposed entropy-guided feature generation method is naturally integrated to the general structure learning framework, thus extending it. We again make use of dependency parsing as an illustrative application to show that the system presented in the previous chapters is an instantiation of the extended framework.

Structure learning is to learn a function that maps an input \mathbf{x} to the correct output structure \mathbf{y} . The output is an arbitrary structure, i.e., it comprises many variables with complex interdependencies. In the SL framework, the prediction problem is recast as an optimization problem of the form

$$F(\mathbf{x}; \mathbf{w}) = \arg \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} s(\mathbf{x}, \mathbf{y}; \mathbf{w}), \quad (4-1)$$

where $\mathbf{w} = (w_1, \dots, w_M)$ is the parameter vector of some learned model, $\mathcal{Y}(\mathbf{x})$ is the set of feasible predictions for the given input, and $s(\mathbf{x}, \mathbf{y}; \mathbf{w})$ is a \mathbf{w} -parameterized scoring function that measures how well an output \mathbf{y} fits the input \mathbf{x} . The prediction output space is arbitrary, but the scoring function must have the following strict form

$$s(\mathbf{x}, \mathbf{y}; \mathbf{w}) = \langle \mathbf{w}, \Phi(\mathbf{x}, \mathbf{y}) \rangle, \quad (4-2)$$

where $\langle \cdot, \cdot \rangle$ is the scalar product operator and $\Phi(\mathbf{x}, \mathbf{y}) = (\phi_1(\mathbf{x}, \mathbf{y}), \dots, \phi_M(\mathbf{x}, \mathbf{y}))$ is some joint feature vector representation of the input-output pair. That is, the scoring function is linear on the feature representation.

4.1

Feature Factorization

Each value $\phi_m(\mathbf{x}, \mathbf{y})$ in the feature vector is called the *global* feature m , which is the value of a specific feature m on the whole structure. For dependency parsing, we define the global feature m as $\phi_m(\mathbf{x}, \mathbf{y}) = \sum_{(i,j) \in \mathbf{y}} \phi_m(\mathbf{x}, i, j)$, which is the sum of the *local* feature values over all *edges* in the tree \mathbf{y} . In that way, we can rewrite the dependency tree scoring function

in the general form of (4-2)

$$\begin{aligned}
s(\mathbf{x}, \mathbf{y}; \mathbf{w}) &= \sum_{(i,j) \in \mathbf{y}} \langle \mathbf{w}, \Phi(\mathbf{x}, i, j) \rangle \\
&= \sum_{(i,j) \in \mathbf{y}} \sum_{m=1}^M w_m \cdot \phi_m(\mathbf{x}, i, j) \\
&= \sum_{m=1}^M w_m \sum_{(i,j) \in \mathbf{y}} \phi_m(\mathbf{x}, i, j) \\
&= \sum_{m=1}^M w_m \cdot \phi_m(\mathbf{x}, \mathbf{y}) \\
&= \langle \mathbf{w}, \Phi(\mathbf{x}, \mathbf{y}) \rangle.
\end{aligned}$$

Thus, global DP features are factored along the edges of the dependency tree. Consequently, the score of a dependency tree in the prediction function is also factored along its edges. Feature factorization is a key point in SL modeling and must give rise to efficient prediction algorithms. For DP, for instance, we end up with a maximum branching problem that is efficiently solved by the Chu-Liu-Edmonds algorithm.

4.2

Entropy-Guided Feature Generation

The derived feature vector $\Phi(\mathbf{x}, \mathbf{y})$ is automatically generated by means of the proposed entropy-guided feature generation method. EFG induces feature templates by conjoining the available basic features and then instantiates these templates to generate the derived feature vectors. Basic features are factored in the same way as derived features. This factorization determines the prediction scoring function and, consequently, directly affects the prediction problem, which is the core of the SL framework. Therefore, we use the same factorization to derive the EFG basic dataset. For each factor in a structured example, we generate an example in the basic dataset, which comprises a vector of basic features and a decision variable. These variables correspond to local decision variables in the prediction problem. Hence, EFG conjoins basic features that help to discriminante the prediction problem variables.

For instance, dependency parsing features are factored along candidate edges. Thus, for each edge in a training sentence, we generate an example in the basic dataset. And, a binary decision variable is associated with each candidate edge and determines whether an edge is present in the corresponding decision

tree. Hence, EFG generates feature templates that are highly discriminative with respect to dependency edge prediction, which corresponds to the local decision in dependency parsing.

Each generated template is used to instantiate a feature in $\Phi(\mathbf{x}, \mathbf{y})$. Then, this derived feature vector is used to train the structured model. Therefore, the structured model is linear on the derived feature representation, which corresponds to a non-linear combination of the basic features. EFG is completely aligned to the SL framework and it is naturally integrated in it.

4.3

Training Algorithm

There are some training algorithms that learn the parameter vector \mathbf{w} from a given training dataset $\mathcal{D} = \{(\mathbf{x}, \mathbf{y})\}$ of correct input-output pairs. For instance, Collins (2002b) proposed the structured perceptron algorithm, a generalization of the well known binary perceptron algorithm for sequence labeling problems. The structured perceptron can be easily applied to any structured problem. Collins (2002b) also proved that the structured perceptron converges to a zero-error solution, if one exists. Crammer and Singer (2003) proposed the margin infused relaxed algorithm (MIRA), an online algorithm to train structured models for multiclass problems. MIRA can also be extended for virtually any structured problem and, for instance, is used in MSTParser. Crammer and Singer (2003) also proved some mistake bounds for an algorithm class called *ultraconservative*, which includes MIRA and structure perceptrons. SVM^{struct} Tsochantaridis et al. (2004) formulates the structure learning problem through a regularized max-margin framework, inspired on the binary support vector machine formulation. They also proposed a cutting plane method to efficiently solve this problem. However, this method still requires more computational power and memory than online algorithms like structure perceptron and MIRA. Additionally, the online algorithms are much simpler to implement than SVM^{struct}.

In this work, we use the structure perceptron algorithm (Collins, 2002b). Given a training sample $\mathcal{D} = \{(\mathbf{x}, \mathbf{y})\}$ of correct input-output pairs, the algorithm generates a sequence of models until convergence. At each iteration, a training instance is drawn from \mathcal{D} and two major steps are performed: prediction using the current model and model update based on the difference between the correct and the predicted outputs. We use the large-margin structure perceptron Fernandes and Brefeld (2011).

During training, instead of the ordinary prediction problem in (4-1), we

use the following loss-augmented version

$$F_\ell(\mathbf{x}; \mathbf{w}) = \arg \max_{\mathbf{y}' \in \mathcal{Y}(\mathbf{x})} [\langle \mathbf{w}, \Phi(\mathbf{x}, \mathbf{y}') \rangle + C \cdot \ell(\mathbf{y}, \mathbf{y}')],$$

where $\ell(\cdot, \cdot) \geq 0$ is an appropriate loss function that measures the difference between a candidate output and the correct one. For dependency trees, as presented in Section 2.2, we use a loss function that just counts the number of incorrect edges in the predicted tree. This loss function factorizes along the dependency tree edges, just like the global features. Thus, the nature of the underlying optimization problem is not modified when using the loss-augmented prediction. This is a desirable, yet *not necessary*, characteristic of loss functions in the SL framework. The model update usually is also factored along the output structure and efficient algorithms can be used. For dependency trees, for instance, this update can be performed in linear time, since a tree has no more than N edges (for a sentence with N tokens plus the artificial token).

We further extend the SL framework by introducing the EFG method into it. EFG is naturally integrated in the SL framework as a preprocessing step. In Figure 4.1, we present the pseudo-code of the entropy-guided large-margin structure perceptron. EFG is used to generate the derived feature vectors

```

 $\Phi(\mathcal{D}) \leftarrow \text{EFG}(\mathcal{D})$ 
 $\mathbf{w}^0 \leftarrow \mathbf{0}$ 
 $t \leftarrow 0$ 
while no convergence
  for each  $(\mathbf{x}, \mathbf{y}) \in \mathcal{D}$ 
     $\hat{\mathbf{y}} \leftarrow \arg \max_{\mathbf{y}' \in \mathcal{Y}(\mathbf{x})} [\langle \mathbf{w}^t, \Phi(\mathbf{x}, \mathbf{y}') \rangle + C \cdot \ell(\mathbf{y}, \mathbf{y}')]$ 
     $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t + \Phi(\mathbf{x}, \mathbf{y}) - \Phi(\mathbf{x}, \hat{\mathbf{y}})$ 
     $t \leftarrow t + 1$ 
return  $\frac{1}{t} \sum_{k=1}^t \mathbf{w}^k$ 

```

Figure 4.1: ESL training algorithm – the entropy-guided large-margin structure perceptron.

$\Phi(\mathcal{D}) = \{\Phi(\mathbf{x}, \mathbf{y})\}_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}}$. The derived features are then used to train the structured model. Note that, when a correct prediction is made, that is $\hat{\mathbf{y}} = \mathbf{y}$, the model does not change, that is $\mathbf{w}^{k+1} \leftarrow \mathbf{w}^k$. When the prediction is wrong, the update rule favors the correct output \mathbf{y} over the predicted one $\hat{\mathbf{y}}$. Regarding binary features, for instance, the update rule increases the weights of features that are present in \mathbf{y} but missing in $\hat{\mathbf{y}}$ and decreases the weights of features that are present in $\hat{\mathbf{y}}$ but not in \mathbf{y} . The weights of features that are present in both \mathbf{y} and $\hat{\mathbf{y}}$ are not changed.

A simple extension of Novikoff's theorem (Novikoff, 1962) shows that the structure perceptron is guaranteed to converge to a zero loss solution, if one exists, in a finite number of steps (Altun et al., 2003; Collins, 2002a). The convergence theorem for SPerc is stated in Theorem 1. Crammer and Singer (2003) further prove some mistake bounds for the structure perceptron algorithm.

Theorem 1 (Structure Perceptron Convergence) *For any training dataset \mathcal{D} that is separable by margin δ , the structure perceptron algorithm makes no more than $\frac{R^2}{\delta^2}$ prediction errors, where $R = \max_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}; \mathbf{y}' \in \mathcal{Y}(\mathbf{x})} \|\Phi(\mathbf{x}, \mathbf{y}) - \Phi(\mathbf{x}, \mathbf{y}')\|$ is the radius of a hypersphere that, centered at $\Phi(\mathbf{x}, \mathbf{y})$, encloses the joint feature vectors for all alternative outputs $\mathbf{y}' \in \mathcal{Y}(\mathbf{x})$, for all training examples $(\mathbf{x}, \mathbf{y}) \in \mathcal{D}$.*

4.4

Kernelization

We argue that EFG has two main advantages over kernel functions. EFG training algorithm is much faster than kernelized algorithms, and EFG makes generalization performance control easier than with kernel methods. In this section, we present the kernelized structure perceptron in order to better understand the differences between this method and EFG.

Analogously to the binary perceptron algorithm, its structure generalization can be easily kernelized. Given the sequence $(\mathbf{x}^1, \mathbf{y}^1, \hat{\mathbf{y}}^1), \dots, (\mathbf{x}^T, \mathbf{y}^T, \hat{\mathbf{y}}^T)$ of inputs, correct outputs and predicted outputs considered by the training algorithm up to iteration T , the parameter vector at this point can be defined as

$$\mathbf{w}^T = \sum_{t=1}^{T-1} [\Phi(\mathbf{x}^t, \mathbf{y}^t) - \Phi(\mathbf{x}^t, \hat{\mathbf{y}}^t)].$$

The algorithm can keep track of how many times each alternative output $\hat{\mathbf{y}}$ has been predicted instead of the correct output \mathbf{y} for each example pair (\mathbf{x}, \mathbf{y}) by means of counters $\alpha_{\mathbf{x}, \mathbf{y}, \hat{\mathbf{y}}}$. Thus, the parameter vector can be rewritten as

$$\mathbf{w} = \sum_{\mathbf{x}, \mathbf{y}, \hat{\mathbf{y}}} \alpha_{\mathbf{x}, \mathbf{y}, \hat{\mathbf{y}}} \cdot [\Phi(\mathbf{x}, \mathbf{y}) - \Phi(\mathbf{x}, \hat{\mathbf{y}})], \quad (4-3)$$

which is called the *dual model representation*. The output space $\mathcal{Y}(\mathbf{x})$ of most SL problems is exponential on the input size or even infinity. Thus, the dual model representation may comprise an intractable number of parameters. However, these parameters are initially *zero* for all $\mathbf{x}, \mathbf{y}, \hat{\mathbf{y}}$ and only need to be instantiated once the respective triple is actually seen.

Using (4-3), the objective function of the prediction problem can also be rewritten as

$$\langle \mathbf{w}, \Phi(\mathbf{x}', \mathbf{y}') \rangle = \sum_{\mathbf{x}, \mathbf{y}, \hat{\mathbf{y}}} \alpha_{\mathbf{x}, \mathbf{y}, \hat{\mathbf{y}}} \cdot [\langle \Phi(\mathbf{x}, \mathbf{y}), \Phi(\mathbf{x}', \mathbf{y}') \rangle - \langle \Phi(\mathbf{x}, \hat{\mathbf{y}}), \Phi(\mathbf{x}', \mathbf{y}') \rangle],$$

which depends only on inner products of feature vectors of the form $\langle \Phi^a, \Phi^b \rangle$, where Φ^a and Φ^b are shortcuts to, respectively, $\Phi(\mathbf{x}^a, \mathbf{y}^a)$ and $\Phi(\mathbf{x}^b, \mathbf{y}^b)$ for any two input-output pairs $(\mathbf{x}^a, \mathbf{y}^a)$ and $(\mathbf{x}^b, \mathbf{y}^b)$. Following the *kernel trick* (Vapnik, 1998), the inner products of feature vectors can then be replaced by an appropriate kernel function

$$K(\Phi^a, \Phi^b) = \langle \Psi(\Phi^a), \Psi(\Phi^b) \rangle,$$

where $\Psi(\cdot)$ expands elements from the original feature vector space $\Phi(\cdot, \cdot)$ to a much higher dimensional space. The kernel trick relies on the kernel function $K(\cdot, \cdot)$ to efficiently compute inner products in the high dimensional space of Ψ without explicitly expanding the original feature space.

The most successful kernel function family for NLP problems is the *polynomial kernel*. Considering binary features, a polynomial kernel of degree d conjoins all possible combinations with up to d original features. The polynomial kernel of degree d can be efficiently computed by

$$K_d(\Phi^a, \Phi^b) = (\langle \Phi^a, \Phi^b \rangle + 1)^d,$$

which involves only an inner product in the original feature space, a sum of a constant, and an exponentiation. For instance, if $d = 2$ and the original space has exactly 3 *binary* features, then the *explicit* polynomial kernel expansion of $\Phi(\mathbf{x}, \mathbf{y}) = (\phi_1, \phi_2, \phi_3)$ corresponds to $\Psi(\Phi(\mathbf{x}, \mathbf{y})) = (1, \phi_1, \phi_2, \phi_3, \phi_1\phi_2, \phi_1\phi_3, \phi_2\phi_3)$, if we omit redundant permutations.

The polynomial kernel of degree d is equivalent to generating all possible templates with length up to d . The problem with these kernels is that the only way to control which combinations are used is through the parameter d . For some SL task, for instance, $d = 2$ can be not enough to capture all relevant contextual patterns, but $d = 3$ can bring so many patterns that it is harmful to the generalization performance. This is a known issue with kernel methods and is related to overfitting. Another issue with kernel functions is training time. Performing predictions with the dual model is much slower than with the primal, because the former is represented by a list of dual variables that usually keeps growing during training. Since the prediction problem is constantly solved during training, the training algorithm becomes very slow.

Just like features, kernel functions can also be decomposed along the output structures. Thus, the dual model representation can be even more sparse by using α counters for each factor that appears in (\mathbf{x}, \mathbf{y}) but not in $(\mathbf{x}, \hat{\mathbf{y}})$, and vice versa. For DP, for instance, we can store a counter for each possible edge within a training sentence.

4.5

Empirical results

We compare ESL to polynomial kernels on two text chunking tasks. We use ESL to train a text chunking system on the Portuguese dataset provided by Fernandes et al. (2010b). We also train a kernelized SPerc system on the same data using a second-degree polynomial kernel. Previous work (Kudo and Matsumoto, 2001; Wu et al., 2006) report that this is the optimum degree for text chunking. Again, in this experiment, we use the same basic features, training algorithm, and datasets for both systems. In the first row of Table 4.1, we report the performances of these two systems. ESL outperforms the kernel

| Task | Kernel Method | | ESL | |
|---------------------|---------------|----------------|----------------|-----------------|
| | Learning | F ₁ | F ₁ | Error Reduction |
| Portuguese Chunking | SPerc | 86.67 | 87.72 | 7.9% |
| English Chunking | SVM | 93.48 | 94.12 | 9.8% |

Table 4.1: Comparison of ESL to second-degree polynomial kernel.

method, reducing its error by 7.9%. That is an impressive achievement, since polynomial kernels present state-of-the-art results on many tasks, and training time for kernel methods is more than one order of magnitude longer than for ESL.

We also compare ESL to a kernel system on the CoNLL-2000 text chunking dataset. The second row of Table 4.1 presents the performances of ESL and a second-degree polynomial kernel system. The kernel method result is reported by Kudo and Matsumoto (2001). They use an SVM algorithm to train their system. However, they employ a training strategy that considers sequential interdependencies among output variables and also use Viterbi decoding during test. ESL outperforms this system, reducing its error by 9.8%.