

# 17 | STRUCTURED PREDICTION

It is often the case that instead of predicting a *single* output, you need to predict multiple, correlated outputs simultaneously. In natural language processing, you might want to assign a syntactic label (like noun, verb, adjective, etc.) to words in a sentence: there is clear correlation among these labels. In computer vision, you might want to label regions in an image with object categories; again, there is correlation among these regions. The branch of machine learning that studies such questions is **structured prediction**.

In this chapter, we will cover two of the most common algorithms for structured prediction: the structured perceptron and the structured support vector machine. We will consider two types of structure. The first is the “sequence labeling” problem, typified by the natural language processing example above, but also common in computational biology (labeling amino acids in DNA) and robotics (labeling actions in a sequence). For this, we will develop specialized prediction algorithms that take advantage of the sequential nature of the task. We will also consider more general structures beyond sequences, and discuss how to cast them in a generic optimization framework: **integer linear programming** (or **ILP**).

The general framework we will explore is that of *jointly scoring input/output configurations*. We will construct algorithms that learn a function  $s(x, \hat{y})$  ( $s$  for “score”), where  $x$  is an input (like an image) and  $\hat{y}$  is some predicted output (like a segmentation of that image). For any given image, there are a *lot* of possible segmentations (i.e., a lot of possible  $\hat{y}$ s), and the goal of  $s$  is to rank them in order of “how good” they are: how compatible they are with the input  $x$ . The most important thing is that the scoring function  $s$  ranks the *true* segmentation  $y$  higher than any other imposter segmentation  $\hat{y}$ . That is, we want to ensure that  $s(x, y) > s(x, \hat{y})$  for all  $\hat{y} \neq y$ . The main challenge we will face is how to do this *efficiently*, given that there are so many imposter  $\hat{y}$ s.

## Learning Objectives:

- Recognize when a problem should be solved using a structured prediction technique.
- Implement the structured perceptron algorithm for sequence labeling.
- Map “argmax” problems to integer linear programs.
- Augment the structured perceptron with losses to derive structured SVMs.

Dependencies:

### 17.1 Multiclass Perceptron

In order to build up to structured problems, let's begin with a simplified by pedagogically useful stepping stone: multiclass classification with a perceptron. As discussed earlier, in multiclass classification we have inputs  $\mathbf{x} \in \mathbb{R}^D$  and output labels  $y \in \{1, 2, \dots, K\}$ . Our goal is to learn a scoring function  $s$  so that  $s(\mathbf{x}, y) > s(\mathbf{x}, \hat{y})$  for all  $\hat{y} \neq y$ , where  $y$  is the true label and  $\hat{y}$  is an imposter label. The general form of scoring function we consider is a linear function of a joint feature vector  $\phi(\mathbf{x}, y)$ :

$$s(\mathbf{x}, y) = \mathbf{w} \cdot \phi(\mathbf{x}, y) \quad (17.1)$$

Here, the features  $\phi(\mathbf{x}, y)$  should denote how "compatible" the input  $\mathbf{x}$  is with the label  $y$ . We keep track of a single weight vector  $\mathbf{w}$  that learns how to weigh these different "compatibility" features.

A natural way to represent  $\phi$ , if you know nothing else about the problem, is an *outer product* between  $\mathbf{x}$  and the label space. This yields the following representation:

$$\phi(\mathbf{x}, k) = \left\langle \underbrace{0, 0, \dots, 0}_{D(k-1) \text{ zeros}}, \underbrace{\mathbf{x}}_{\in \mathbb{R}^D}, \underbrace{0, 0, \dots, 0}_{D(K-k) \text{ zeros}} \right\rangle \in \mathbb{R}^{DK} \quad (17.2)$$

In this representation,  $\mathbf{w}$  effectively encodes a separate weight for every feature/label pair.

How are we going to learn  $\mathbf{w}$ ? We will start with  $\mathbf{w} = \mathbf{0}$  and then process each input one at a time. Suppose we get an input  $\mathbf{x}$  with gold standard label  $y$ . We will use the current scoring function to predict a label. In particular, we will predict the label  $\hat{y}$  that maximizes the score:

$$\hat{y} = \operatorname{argmax}_{\hat{y} \in [1, K]} s(\mathbf{x}, \hat{y}) \quad (17.3)$$

$$= \operatorname{argmax}_{\hat{y} \in [1, K]} \mathbf{w} \cdot \phi(\mathbf{x}, \hat{y}) \quad (17.4)$$

If this predicted output is correct (i.e.,  $\hat{y} = y$ ), then, per the normal perceptron, we will do nothing. Suppose that  $\hat{y} \neq y$ . This means that the score of  $\hat{y}$  is greater than the score of  $y$ , so we want to update  $\mathbf{w}$  so that the score of  $\hat{y}$  is decreased and the score of  $y$  is increased. We do this by:

$$\mathbf{w} \leftarrow \mathbf{w} + \phi(\mathbf{x}, y) - \phi(\mathbf{x}, \hat{y}) \quad (17.5)$$

To make sure this is doing what we expect, let's consider what would happen if we computed scores under the *updated* value of  $\mathbf{w}$ . To make the notation clear, let's say  $\mathbf{w}^{(\text{old})}$  are the weights before update, and

**Algorithm 39** MULTICLASSPERCEPTRONTRAIN( $\mathbf{D}$ ,  $MaxIter$ )

---

```

1:  $w \leftarrow \mathbf{0}$  // initialize weights
2: for  $iter = 1 \dots MaxIter$  do
3:   for all  $(x, y) \in \mathbf{D}$  do
4:      $\hat{y} \leftarrow \operatorname{argmax}_k w \cdot \phi(x, k)$  // compute prediction
5:     if  $\hat{y} \neq y$  then
6:        $w \leftarrow w + \phi(x, y) - \phi(x, \hat{y})$  // update weights
7:     end if
8:   end for
9: end for
10: return  $w$  // return learned weights

```

---

$w^{(\text{new})}$  are the weights after update. Then:

$$w^{(\text{new})} \cdot \phi(x, y) \quad (17.6)$$

$$= \left( w^{(\text{old})} + \phi(x, y) - \phi(x, \hat{y}) \right) \cdot \phi(x, y) \quad (17.7)$$

$$= \underbrace{w^{(\text{old})} \cdot \phi(x, y)}_{\text{old prediction}} + \underbrace{\phi(x, y) \cdot \phi(x, y)}_{\geq 0} - \underbrace{\phi(x, \hat{y}) \cdot \phi(x, y)}_{=0} \quad (17.8)$$

Here, the first term is the old prediction. The second term is of the form  $a \cdot a$  which is non-negative (and, unless  $\phi(x, y)$  is the zero vector, positive). The third term is the dot product between  $\phi$  for two different labels, which by definition of  $\phi$  is zero (see Eq (17.2)).

This gives rise to the updated multiclass perceptron specified in Algorithm 17.1. As with the normal perceptron, the generalization of the multiclass perceptron increases dramatically if you do weight averaging.

An important note is that MulticlassPerceptronTrain is actually more powerful than suggested so far. For instance, suppose that you have three categories, but believe that two of them are tightly related, while the third is very different. For instance, the categories might be {music, movies, oncology}. You can encode this relatedness by defining a feature expansion  $\phi$  that reflects this:

$$\phi(x, \text{music}) = \langle x, \mathbf{0}, \mathbf{0}, x \rangle \quad (17.9)$$

$$\phi(x, \text{movies}) = \langle \mathbf{0}, x, \mathbf{0}, x \rangle \quad (17.10)$$

$$\phi(x, \text{oncology}) = \langle \mathbf{0}, \mathbf{0}, x, \mathbf{0} \rangle \quad (17.11)$$

This encoding is identical to the normal encoding in the first three positions, but includes an extra copy of the features at the end, shared between music and movies. By doing so, if the perceptron wants to learn something common to music and movies, it can use this final shared position.

**?** Verify the score of  $\hat{y}$ ,  $w^{(\text{new})} \cdot \phi(x, \hat{y})$ , decreases after an update, as we would want.

**?** Suppose you have a *hierarchy* of classes arranged in a tree. How could you use that to construct a feature representation. You can think of the music/movies/oncology example as a binary tree: the left branch of the root splits into music and movies; the right branch of the root is just oncology.

## 17.2 Structured Perceptron

Let us now consider the sequence labeling task. In sequence labeling, the outputs are themselves variable-length vectors. An input/output pair (which must have the same length) might look like:

$$\mathbf{x} = \text{“ monsters eat tasty bunnies “} \quad (17.12)$$

$$\mathbf{y} = \text{noun verb adj noun} \quad (17.13)$$

To set terminology, we will refer to the *entire sequence*  $\mathbf{y}$  as the “output” and a single label within  $\mathbf{y}$  as a “label”. As before, our goal is to learn a scoring function that scores the true output sequence  $\mathbf{y}$  higher than any imposter output sequence.

As before, despite the fact that  $\mathbf{y}$  is now a vector, we can *still* define feature functions over the *entire* input/output pair. For instance, we might want to count the number of times “monsters” has been tagged as “noun” in a given output. Or the number of times “verb” is followed by “noun” in an output. Both of these are features that are likely indicative of a *correct* output. We might also count the number of times “tasty” has been tagged as a verb (probably a negative feature) and the number of times two verbs are adjacent (again, probably a negative feature).

More generally, a very standard set of features would be:

- the number of times word  $w$  has been labeled with tag  $l$ , for all words  $w$  and all syntactic tags  $l$
- the number of times tag  $l$  is adjacent to tag  $l'$  in the output, for all tags  $l$  and  $l'$

The first set of features are often called **unary features**, because they talk only about the relationship between the input (sentence) and a *single* (unit) label in the output sequence. The second set of features are often called **Markov features**, because they talk about adjacent labels in the output sequence, which is reminiscent of Markov models which only have short term memory.

Note that for a given input  $x$  of length  $L$  (in the example,  $L = 4$ ), the number of possible outputs is  $K^L$ , where  $K$  is the number of syntactic tags. This means that the number of possible outputs *grows exponentially* in the length of the input. In general, we write  $\mathcal{Y}(x)$  to mean “the set of all possible structured outputs for the input  $x$ ”. We have just seen that  $|\mathcal{Y}(x)| = K^{\text{len}(x)}$ .

Despite the fact that the inputs and outputs have variable length, the size of the *feature representation* is constant. If there are  $V$  words in your vocabulary and  $K$  labels for a given word, the the number of unary features is  $VK$  and the number of Markov features is  $K^2$ , so

**Algorithm 40** STRUCTUREDPERCEPTRONTRAIN( $\mathbf{D}$ ,  $MaxIter$ )

---

```

1:  $w \leftarrow \mathbf{0}$  // initialize weights
2: for  $iter = 1 \dots MaxIter$  do
3:   for all  $(x, y) \in \mathbf{D}$  do
4:      $\hat{y} \leftarrow \operatorname{argmax}_{\hat{y} \in \mathcal{Y}(x)} w \cdot \phi(x, \hat{y})$  // compute prediction
5:     if  $\hat{y} \neq y$  then
6:        $w \leftarrow w + \phi(x, y) - \phi(x, \hat{y})$  // update weights
7:     end if
8:   end for
9: end for
10: return  $w$  // return learned weights

```

---

the total number of features is  $K(V + K)$ . Of course, more complex feature representations are possible and, in general, are a good idea. For example, it is often useful to have unary features of neighboring words like “the number of times the word immediately preceding a verb was ‘monsters.’”

Now that we have a fixed size feature representation, we can develop a perceptron-style algorithm for sequence labeling. The core idea is the same as before. We will maintain a *single* weight vector  $w$ . We will make predictions by choosing the (entire) output sequence  $\hat{y}$  that maximizes a score given by  $w \cdot \phi(x, \hat{y})$ . And if this output sequence is incorrect, we will adjust the weights *word* the correct output sequence  $y$  and away from the incorrect output sequence  $\hat{y}$ . This is summarized in Algorithm 17.2

You may have noticed that Algorithm 17.2 for the structured perceptron is *identical* to Algorithm 17.1, aside from the fact that in the multiclass perceptron the argmax is over the  $K$  possible classes, while in the structured perceptron, the argmax is over the  $K^L$  possible output sequences!

The only difficulty in this algorithm is in line 4:

$$\hat{y} \leftarrow \operatorname{argmax}_{\hat{y} \in \mathcal{Y}(x)} w \cdot \phi(x, \hat{y}) \quad (17.14)$$

In principle, this requires you to search over  $K^L$  possible output sequences  $\hat{y}$  to find the one that maximizes the dot product. Except for very small  $K$  or very small  $L$ , this is computationally infeasible. Because of its difficulty, this is often referred to as the **argmax problem** in structured prediction. Below, we consider how to solve the argmax problem for sequences.

### 17.3 Argmax for Sequences

We now face an *algorithmic* question, not a machine learning question: how to compute the argmax in Eq 17.14 efficiently. In general,

this is not possible. However, under somewhat restrictive assumptions about the form of our features  $\phi$ , we can solve this problem efficiently, by casting it as the problem of computing a maximum weight path through a specifically constructed lattice. This is a variant of the Viterbi algorithm for hidden Markov models, a classic example of dynamic programming. (Later, in Section 17.6, we will consider argmax for more general problems.)

The key observation for sequences is that—so long as we restrict our attention to unary features and Markov features—the feature function  $\phi$  decomposes over the input. This is easiest to see with an example. Consider the input/output sequence from before:  $x =$  “monsters eat tasty bunnies” and  $y =$  [noun verb adj noun]. If we want to compute the number of times “bunnies” is tagged as “noun” in this pair, we can do this by:

1. count the number of times “bunnies” is tagged as “noun” in the first three words of the sentence
2. add to that the number of times “bunnies” is tagged as “noun” in the final word

We can do a similar exercise for Markov features, like the number of times “adj” is followed by “noun”.

However, we don’t actually *need* these counts. All we need for computing the argmax sequence is the dot product between the weights  $w$  and these counts. In particular, we can compute  $w \cdot \phi(x, y)$  as the dot product on all-but-the-last word plus the dot product on the last word:  $w \cdot \phi_{1:3}(x, y) + w \cdot \phi_4(x, y)$ . Here,  $\phi_{1:3}$  means “features for everything up to and including position 3” and  $\phi_4$  means “features for position 4.”

More generally, we can write  $\phi(x, y) = \sum_{l=1}^L \phi_l(x, y)$ , where  $\phi_l(x, y)$  only includes features about position  $l$ .<sup>1</sup> In particular, we’re taking advantage of the associative law for addition:

$$w \cdot \phi(x, y) = w \cdot \sum_{l=1}^L \phi_l(x, y) \quad \text{decomposition of structure} \quad (17.15)$$

$$= \sum_{l=1}^L w \cdot \phi_l(x, y) \quad \text{associative law} \quad (17.16)$$

What this means is that we can build a graph like that in Figure ??, with one verticle slice per time step ( $l = 1 \dots L$ ).<sup>2</sup> Each *edge* in this graph will receive a weight, constructed in such a way that if you take a complete path through the lattice, and add up all the weights, this will correspond exactly to  $w \cdot \phi(x, y)$ .

To complete the construction, let  $\phi_l(x, \dots \circ y \circ y')$  denote the *unary* features at position  $l$  together with the *Markov* features that end at

<sup>1</sup> In the case of Markov features, we think of them as pairs that *end* at position  $l$ , so “verb adj” would be the active feature for  $\phi_3$ .

<sup>2</sup> A graph of this sort is called a **trellis**, and sometimes a **lattice** in the literature.

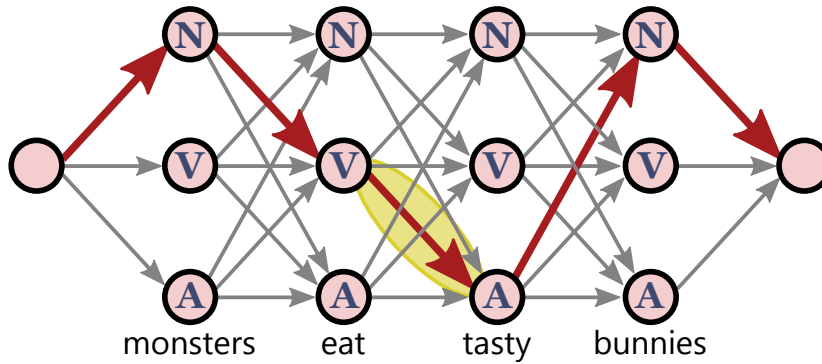


Figure 17.1: A picture of a trellis sequence labeling. At each time step  $l$  the corresponding word can have any of the three possible labels. Any path through this trellis corresponds to a unique labeling of this sentence. The gold standard path is drawn with bold red arrows. The highlighted edge corresponds to the edge between  $l = 2$  and  $l = 3$  for verb/adj as described in the text. That edge has weight  $w \cdot \phi_3(x, \dots \circ \text{verb} \circ \text{adj})$ .

position  $l$ . These features depend *only* on  $x$ ,  $y$  and  $y'$ , and *not* any of the previous parts of the output.

For example, in the running example “monsters/noun eat/verb tasty/adj bunnies/noun”, consider the edge between  $l = 2$  and  $l = 3$  going from “verb” to “adj”. (Note: this is a “correct” edge, in the sense that it belongs to the ground truth output.) The features associated with this edge will be unary features about “tasty/adj” as well as Markov features about “verb/adj”. The *weight* of this edge will be exactly the total score (according to  $w$ ) of those features.

Formally, consider an edge in the trellis that goes from time  $l - 1$  to  $l$ , and transitions from  $y$  to  $y'$ . Set the weight of this edge to exactly  $w \cdot \phi_l(x, \dots \circ y \circ y')$ . By doing so, we guarantee that the sum of weights along any path through this lattice is exactly equal to the score of that path. Once we have constructed the graph as such, we can run any max-weight path algorithm to compute the highest scoring output. For trellises, this can be computed by the Viterbi algorithm, or by applying any of a number of path finding algorithms for more general graphs. A complete derivation of the dynamic program in this case is given in Section 17.7 for those who want to implement it directly.

The main benefit of this construction is that it is guaranteed to exactly compute the argmax output for sequences required in the structured perceptron algorithm, *efficiently*. In particular, it’s runtime is  $O(LK^2)$ , which is an exponential improvement on the naive  $O(K^L)$  runtime if one were to enumerate every possible output sequence. The algorithm can be naturally extended to handle “higher order” Markov assumptions, where features depend on triples or quadruples of the output. The trellis becomes larger, but the algorithm remains essentially the same. In order to handle a length  $M$  Markov features, the resulting algorithm will take  $O(LK^M)$  time. In practice, it’s rare that  $M > 3$  is necessary or useful.

## 17.4 Structured Support Vector Machines

In Section 8.7 we saw the support vector machine as a very useful general framework for binary classification. In this section, we will develop a related framework for structured support vector machines. The two main advantages of structured SVMs over the structured perceptron are (1) it is regularized (though averaging in structured perceptron achieves a similar effect) and (2) we can incorporate more complex loss functions.

In particular, one suboptimal thing about the structured perceptron is that all errors are considered equally bad. For structured problems, we often have much more nuanced and elaborate loss functions that we want to optimize. Even for sequence labeling, it is typically far worse to label every word incorrectly than to just label one word incorrectly. It is very common to use **Hamming loss** as a general loss function for structured prediction. Hamming loss simply counts: of all the predictions you made, how many were incorrect? For sequence labeling, it is:

$$\ell^{(\text{Ham})}(\mathbf{y}, \hat{\mathbf{y}}) = \sum_{l=1}^L \mathbf{1}[y_l \neq \hat{y}_l] \quad (17.17)$$

In order to build up to structured SVMs, recall that SVMs began with the following optimization problem:

$$\begin{aligned} \min_{\mathbf{w}, \xi} \quad & \underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{large margin}} + C \underbrace{\sum_n \xi_n}_{\text{small slack}} & (17.18) \\ \text{subj. to} \quad & y_n (\mathbf{w} \cdot \mathbf{x}_n + b) \geq 1 - \xi_n & (\forall n) \\ & \xi_n \geq 0 & (\forall n) \end{aligned}$$

After a bit of work, we were able to reformulate this in terms of a standard loss optimization algorithm with hinge loss:

$$\min_{\mathbf{w}} \quad \underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{large margin}} + C \underbrace{\sum_n \ell^{(\text{hin})}(y_n, \mathbf{w} \cdot \mathbf{x}_n + b)}_{\text{small slack}} \quad (17.19)$$

We can do a similar derivation in the structured case. The question is: exactly what should the slack be measuring? Our *goal* is for the score of the true output  $\mathbf{y}$  to beat the score of any imposter output  $\hat{\mathbf{y}}$ . To incorporate loss, we will say that we want the score of the true output to beat the score of any imposter output by *at least* the loss that would be suffered if we were to predict that imposter output. An alternative view is the ranking view: we want the true output to be ranked above any imposter by an amount at least equal to the loss.



To keep notation simple, we will write  $s_w(\mathbf{x}, \mathbf{y})$  to denote the score of the pair  $\mathbf{x}, \mathbf{y}$ , namely  $\mathbf{w} \cdot \phi(\mathbf{x}, \mathbf{y})$ . This suggests a set of constraints of the form:

$$s_w(\mathbf{x}, \mathbf{y}) - s_w(\mathbf{x}, \hat{\mathbf{y}}) \geq \ell^{(\text{Ham})}(\mathbf{y}, \hat{\mathbf{y}}) - \zeta_{\hat{\mathbf{y}}} \quad (\forall n, \forall \hat{\mathbf{y}} \in \mathcal{Y}(\mathbf{x})) \quad (17.20)$$

The rest of the optimization problem remains the same, yielding:

$$\min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_n \sum_{\hat{\mathbf{y}} \in \mathcal{Y}(\mathbf{x}_n)} \zeta_{n, \hat{\mathbf{y}}} \quad (17.21)$$

$$\begin{aligned} \text{subj. to } & s_w(\mathbf{x}, \mathbf{y}) - s_w(\mathbf{x}, \hat{\mathbf{y}}) \\ & \geq \ell^{(\text{Ham})}(\mathbf{y}_n, \hat{\mathbf{y}}) - \zeta_{n, \hat{\mathbf{y}}} \quad (\forall n, \forall \hat{\mathbf{y}} \in \mathcal{Y}(\mathbf{x}_n)) \\ & \zeta_{n, \hat{\mathbf{y}}} \geq 0 \quad (\forall n, \forall \hat{\mathbf{y}} \in \mathcal{Y}(\mathbf{x}_n)) \end{aligned}$$

This optimization problem asks for a large margin and small slack, where there is a slack very for every training example and every possible incorrect output associated with that training example. In general, this is *way too many* slack variables and *way too many* constraints!

There is a very useful, general trick we can apply. If you focus on the first constraint, it roughly says (letting  $s(\cdot)$  denote score):  $s(\mathbf{y}) \geq [s(\hat{\mathbf{y}}) + \ell(\mathbf{y}, \hat{\mathbf{y}})]$  for all  $\hat{\mathbf{y}}$ , modulo slack. We'll refer to the thing in brackets as the "loss-augmented score." But if we want to guarantee that the score of the true  $\mathbf{y}$  beats the loss-augmented score of *all*  $\hat{\mathbf{y}}$ , it's enough to ensure that it beats the loss-augmented score of the most confusing imposter. Namely, it is sufficient to require that  $s(\mathbf{y}) \geq \max_{\hat{\mathbf{y}}} [s(\hat{\mathbf{y}}) + \ell(\mathbf{y}, \hat{\mathbf{y}})]$ , modulo slack. Expanding out the definition of  $s(\cdot)$  and adding slack back in, we can replace the exponentially large number of constraints in Eq (17.21) with the simpler set of constraints:

$$s_w(\mathbf{x}_n, \mathbf{y}_n) \geq \max_{\hat{\mathbf{y}} \in \mathcal{Y}(\mathbf{x}_n)} [s_w(\mathbf{x}_n, \hat{\mathbf{y}}) + \ell^{(\text{Ham})}(\mathbf{y}_n, \hat{\mathbf{y}})] - \zeta_n \quad (\forall n)$$

We can now apply the same trick as before to remove  $\zeta_n$  from the analysis. In particular, because  $\zeta_n$  is constrained to be  $\geq 0$  and because we are trying to minimize it's sum, we can figure out that out the optimum, it will be the case that:

$$\zeta_n = \max \left\{ 0, \max_{\hat{\mathbf{y}} \in \mathcal{Y}(\mathbf{x}_n)} [s_w(\mathbf{x}_n, \hat{\mathbf{y}}) + \ell^{(\text{Ham})}(\mathbf{y}_n, \hat{\mathbf{y}})] - s_w(\mathbf{x}_n, \mathbf{y}_n) \right\} \quad (17.22)$$

$$= \ell^{(\text{s-h})}(\mathbf{y}_n, \mathbf{x}_n, \mathbf{w}) \quad (17.23)$$

This value is referred to as the **structured hinge loss**, which we have denoted as  $\ell^{(\text{s-h})}(\mathbf{y}_n, \mathbf{x}_n, \mathbf{w})$ . This is because, although it is more complex, it bears a striking resemblance to the **hinge loss** from Chapter 8.

In particular, if the score of the true output beats the score of every the best imposter by at least its loss, then  $\zeta_n$  will be zero. On the other hand, if some imposter (plus its loss) beats the true output, the loss scales linearly as a function of the difference. At this point, there is nothing special about Hamming loss, so we will replace it with some arbitrary structured loss  $\ell$ .

Plugging this back into the objective function of Eq (17.21), we can write the structured SVM as an *unconstrained* optimization problem, akin to Eq (17.19), as:

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_n \ell^{(s-h)}(\mathbf{y}_n, \mathbf{x}_n, \mathbf{w}) \quad (17.24)$$

This is now in a form that we can optimize using subgradient descent (Chapter 8) or stochastic subgradient descent (Chapter 14).

In order to compute subgradients of Eq (17.24), we need to be able to compute subgradients of the structured hinge loss. Mathematically this is straightforward. If the structured hinge loss on an example  $(\mathbf{x}, y)$  is zero, then the gradient with respect to  $\mathbf{w}$  is also zero. If the structured hinge loss is positive, then the gradient is:

$$\nabla_{\mathbf{w}} \ell^{(s-h)}(\mathbf{y}, \mathbf{x}, \mathbf{w}) \quad \text{if the loss is } > 0 \quad (17.25)$$

expand definition using arbitrary structured loss  $\ell$

$$= \nabla_{\mathbf{w}} \left\{ \max_{\hat{\mathbf{y}} \in \mathcal{Y}(\mathbf{x}_n)} [\mathbf{w} \cdot \phi(\mathbf{x}_n, \hat{\mathbf{y}}) + \ell(\mathbf{y}_n, \hat{\mathbf{y}})] - \mathbf{w} \cdot \phi(\mathbf{x}_n, \mathbf{y}_n) \right\} \quad (17.26)$$

define  $\hat{\mathbf{y}}_n$  to be the output that attains the maximum above, rearrange

$$= \nabla_{\mathbf{w}} \left\{ \mathbf{w} \cdot \phi(\mathbf{x}_n, \hat{\mathbf{y}}) - \mathbf{w} \cdot \phi(\mathbf{x}_n, \mathbf{y}_n) + \ell(\mathbf{y}_n, \hat{\mathbf{y}}) \right\} \quad (17.27)$$

take gradient

$$= \phi(\mathbf{x}_n, \hat{\mathbf{y}}) - \phi(\mathbf{x}_n, \mathbf{y}_n) \quad (17.28)$$

Putting this together, we get the full gradient as:

$$\nabla_{\mathbf{w}} \ell^{(s-h)}(\mathbf{y}_n, \mathbf{x}_n, \mathbf{w}) = \begin{cases} \mathbf{0} & \text{if } \ell^{(s-h)}(\mathbf{y}_n, \mathbf{x}_n, \mathbf{w}) = 0 \\ \phi(\mathbf{x}_n, \hat{\mathbf{y}}_n) - \phi(\mathbf{x}_n, \mathbf{y}_n) & \text{otherwise} \end{cases} \quad (17.29)$$

where  $\hat{\mathbf{y}}_n = \operatorname{argmax}_{\hat{\mathbf{y}}_n \in \mathcal{Y}(\mathbf{x}_n)} [\mathbf{w} \cdot \phi(\mathbf{x}_n, \hat{\mathbf{y}}_n) + \ell(\mathbf{y}_n, \hat{\mathbf{y}}_n)]$

The form of this gradient is very simple: it is equal to the features of the worst imposter minus the features of the truth, unless the truth beats all imposters, in which case it's zero. When plugged into stochastic subgradient descent, you end up with an update that looks very much like the structured perceptron: if the current prediction ( $\hat{\mathbf{y}}_n$ ) is correct, there is no gradient step. But if the current prediction is incorrect, you step  $\mathbf{w}$  toward the truth and away from the imposter.

**Algorithm 41** STOCHSUBGRADSTRUCTSVM( $\mathbf{D}$ ,  $MaxIter$ ,  $\lambda$ ,  $\ell$ )

---

```

1:  $w \leftarrow \mathbf{0}$  // initialize weights
2: for  $iter = 1 \dots MaxIter$  do
3:   for all  $(x, y) \in \mathbf{D}$  do
4:      $\hat{y} \leftarrow \operatorname{argmax}_{\hat{y} \in \mathcal{Y}(x)} w \cdot \phi(x, \hat{y}) + \ell(y, \hat{y})$  // loss-augmented prediction
5:     if  $\hat{y} \neq y$  then
6:        $w \leftarrow w + \phi(x, y) - \phi(x, \hat{y})$  // update weights
7:     end if
8:      $w \leftarrow w - \frac{\lambda}{N} w$  // shrink weights due to regularizer
9:   end for
10: end for
11: return  $w$  // return learned weights

```

---

We will consider how to compute the loss-augmented argmax in the next section, but before that we summarize an algorithm for optimizing structured SVMs using stochastic subgradient descent: Algorithm 17.4. Of course there are other possible optimization strategies; we are highlighting this one because it is nearly identical to the structured perceptron. The only differences are: (1) on line 4 you use loss-augmented argmax instead of argmax; and (2) on line 8 the weights are shrunk slightly corresponding to the  $\ell_2$  regularizer on  $w$ . (Note: we have used  $\lambda = 1/(2C)$  to make the connection to linear models clearer.)

### 17.5 Loss-Augmented Argmax

The challenge that arises is that we now have a more complicated argmax problem than before. In structured perceptron, we only needed to compute  $\hat{y}_n$  as the output that maximized its score (see Eq 17.14). Here, we need to find the output that maximizes its score *plus* its loss (Eq (17.29)). This optimization problem is referred to as **loss-augmented search** or **loss-augmented inference**.

Before solving the loss-augmented inference problem, it's worth thinking about why it makes sense. What is  $\hat{y}_n$ ? It's the output that has the highest score among all outputs, *after* adding the output's corresponding loss to that score. In other words, every incorrect output gets an artificial boost to its score, equal to its loss. The loss is serving to make imposters look *even better* than they really are, so if the truth is to beat an imposter, it has to beat it by a *lot*. In fact, this loss augmentation is essentially playing the role of a margin, where the required margin scales according to the loss.

The algorithmic question, then, is how to compute  $\hat{y}_n$ . In the fully general case, this is at least as hard as the normal argmax problem, so we cannot expect a general solution. Moreover, even in cases where the argmax problem is easy (like for sequences), the loss-augmented

argmax problem can still be difficult. In order to make it easier, we need to assume that the loss *decomposes* of the input in a way that's consistent with the features. In particular, if the structured loss function is Hamming loss, this is often straightforward.

As a concrete example, let's consider loss-augmented argmax for sequences under Hamming loss. In comparison to the trellis problem solved in Section 17.7, the only difference is that we want to *reward* paths that go through incorrect nodes in the trellis! In particular, in Figure 17.1, all of the edges that are not part of the gold standard path—those that are thinner and grey—get a free “+1” added to their weights. Since Hamming loss adds one to the score for any word that's predicted incorrectly, this means that every edge in the trellis that leads to an *incorrect* node (i.e., one that does not match the gold truth label) gets a “+1” added to its weight.

Again, consider an edge in the trellis that goes from time  $l - 1$  to  $l$ , and transitions from  $y$  to  $y'$ . In the non-loss-augmented, the weight of this edge was exactly  $w \cdot \phi_l(x, \dots \circ y \circ y')$ . In the loss-augmented cases, the weight of this edge becomes:

$$\underbrace{w \cdot \phi_l(x, \dots \circ y \circ y')}_{\text{edge score, as before}} + \underbrace{\mathbf{1}[y' \neq y_l]}_{+1 \text{ for mispredictions}} \quad (17.30)$$

Once this loss-augmented graph has been constructed, the same max-weight path algorithm can be run to find the loss-augmented argmax sequence.

## 17.6 Argmax in General

The general argmax problem for structured perceptron is the algorithmic question of whether the following can be efficiently computed:

$$\hat{y} \leftarrow \operatorname{argmax}_{\hat{y} \in \mathcal{Y}(x)} w \cdot \phi(x, \hat{y}) \quad (17.31)$$

We have seen that *if* the output space  $\mathcal{Y}(x)$  is sequences *and* the only types of features are unary features and Markov features, then this can be computed efficiently. There are a small number of other structured output spaces and feature restrictions for which efficient problem-specific algorithms exist:

- Binary trees, with context-free features: use the CKY algorithm
- 2d image segmentation, with adjacent-pixel features: use a form of graph cuts
- Spanning trees, with edge-based features: use Kruskal's algorithm (or for directed spanning trees, use Chu-Liu/Edmonds algorithm)

These special cases are often very useful, and many problems can be cast in one of these frameworks. However, it is often the case that you need a more general solution.

One of the most generally useful solutions is to cast the argmax problem as an **integer linear program**, or **ILP**. ILPs are a specific type of mathematical program/optimization problem, in which the objective function being optimized is linear and the constraints are linear. However, unlike “normal” linear programs, in an ILP you are allowed to have integer constraints and disallow fractional values. The general form of an ILP is, for a fixed vector  $a$ :

$$\max_z \quad a \cdot z \quad \text{subj. to} \quad \text{linear constraints on } z \quad (17.32)$$

The main point is that the constraints on  $z$  are allowed to include constraints like  $z_3 \in \{0, 1\}$ , which is considered an integer constraint.

Being able to cast your argmax problem as an ILP has the advantage that there are very good, efficiently, well-engineered ILP solvers out there in the world.<sup>3</sup> ILPs are not a panacea though: in the worst case, the ILP solver will be horribly inefficient. But for prototyping, or if there are no better options, it’s a very handy technique.

<sup>3</sup> I like Gurobi best, and it’s free for academic use. It also has a really nice Python interface.

Figuring out how exactly to cast your argmax problem as an ILP can be a bit challenging. Let’s start with an example of encoding sequence labeling with Markov features as an ILP. We first need to decide what the variables will be. Because we need to encode pairwise features, we will let our variables be of the form:

$$z_{l,k',k} = \mathbf{1}[\text{label } l \text{ is } k \text{ and label } l-1 \text{ is } k'] \quad (17.33)$$

These  $z$ s will all be binary indicator variables.

Our next task is to construct the linear objective function. To do so, we need to assign a value to  $a_{l,k',k}$  in such a way that  $a \cdot z$  will be exactly equal to  $w \cdot \phi(x, y(z))$ , where  $y(z)$  denotes the sequence that we can read off of the variables  $z$ . With a little thought, we arrive at:

$$a_{l,k',k} = w \cdot \phi_l(x, \langle \dots, k', k \rangle) \quad (17.34)$$

Finally, we need to construct constraints. There are a few things that these constraints need to enforce:

1. That all the  $z$ s are binary. That’s easy: just say  $z_{l,k',k} \in \{0, 1\}$ , for all  $l, k', k$ .
2. That for a given position  $l$ , there is exactly one active  $z$ . We can do this with an equality constraint:  $\sum_k \sum_{k'} z_{l,k',k} = 1$  for all  $l$ .
3. That the  $z$ s are internally consistent: if the label at position 5 is supposed to be “noun” then both  $z_{5,..}$  and  $z_{6,..}$  need to agree on

this. We can do this as:  $\sum_{k'} z_{l,k',k} = \sum_{k''} z_{l+1,k,k''}$  for all  $l, k$ . Effectively what this is saying is that  $z_{5,?,verb} = z_{6,verb,?}$  where the “?” means “sum over all possibilities.”

This fully specifies an ILP that you can relatively easily implement (arguably more easily than the dynamic program in Algorithm 17.7) and which will solve the argmax problem for you. Will it be efficient? In this case, probably yes. Will it be as efficient as the dynamic program? Probably not.

It takes a bit of effort and time to get used to casting optimization problems as ILPs, and certainly not all can be, but most can and it's a very nice alternative.

In the case of loss-augmented search for structured SVMs (as opposed to structured perceptron), the objective function of the ILP will need to be modified to include terms corresponding to the loss.

## 17.7 Dynamic Programming for Sequences

Recall the decomposition we derived earlier:

$$w \cdot \phi(x, y) = w \cdot \sum_{l=1}^L \phi_l(x, y) \quad \text{decomposition of structure} \quad (17.35)$$

$$= \sum_{l=1}^L w \cdot \phi_l(x, y) \quad \text{associative law} \quad (17.36)$$

This decomposition allows us to construct the following dynamic program. We will compute  $\alpha_{l,k}$  as the score of the *best possible* output prefix up to and including position  $l$  that labels the  $l$ th word with label  $k$ . More formally:

$$\alpha_{l,k} = \max_{\hat{y}_{1:l-1}} w \cdot \phi_{1:l}(x, \hat{y} \circ k) \quad (17.37)$$

Here,  $\hat{y}$  is a sequence of length  $l - 1$ , and  $\hat{y} \circ k$  denotes the sequence of length  $l$  obtained by adding  $k$  onto the end. The max denotes the fact that we are seeking the *best possible* prefix up to position  $l - 1$ , and the forcing the label for position  $l$  to be  $k$ .

Before working through the details, let's consider an example. Suppose that we've computing the  $\alpha$ s up to  $l = 2$ , and have:  $\alpha_{2,noun} = 2$ ,  $\alpha_{2,verb} = 9$ ,  $\alpha_{2,adj} = -1$  (recall: position  $l = 2$  is “eat”). We want to extend this to position 3; for example, we want to compute  $\alpha_{3,adj}$ . Let's assume there's a single unary feature here, “tasty/adj” and three possible Markov features of the form “?:adj”. Assume these weights are as given to the right.<sup>4</sup> Now, the question for  $\alpha_{3,adj}$  is: what's the score of the best prefix that labels “tasty” as “adj”? We can obtain this by taking the best prefix up to “eat” and then appending

<sup>4</sup>  $w^{\text{“tasty/adj”}} = 1.2$   
 $w^{\text{“noun:adj”}} = -5$   
 $w^{\text{“verb:adj”}} = 2.5$   
 $w^{\text{“adj:adj”}} = 2.2$

each possible label. Whichever combination is best is the winner. The relevant computation is:

$$\alpha_{3,\text{adj}} = \max \left\{ \begin{aligned} &\alpha_{2,\text{noun}} + w_{\text{tasty/adj}} + w_{\text{noun:adj}} \\ &\alpha_{2,\text{verb}} + w_{\text{tasty/adj}} + w_{\text{verb:adj}} \\ &\alpha_{2,\text{adj}} + w_{\text{tasty/adj}} + w_{\text{adj:adj}} \end{aligned} \right\} \quad (17.38)$$

$$= \max \left\{ 2 + 1.2 - 5, \quad 9 + 1.2 + 2.5, \quad -1 + 1.2 + 2.2 \right\} \quad (17.39)$$

$$= \max \left\{ -1.8, \quad 12.7, \quad 2.4 \right\} = 12.7 \quad (17.40)$$

This means that (a) the score for the prefix ending at position 3 labeled as adjective is 12.7, and (b) the “winning” previous label was “verb”. We will need to record these winning previous labels so that we can extract the best path at the end. Let’s denote by  $\zeta_{l,k}$  the label at position  $l - 1$  that achieves the max.

From here, we can formally compute the  $\alpha$ s recursively. The main observation that will be necessary is that, because we have limited ourselves to Markov features,  $\phi_{l+1}(x, \langle y_1, y_2, \dots, y_l, y_{l+1} \rangle)$  depends only on the last two terms of  $y$ , and does *not* depend on  $y_1, y_2, \dots, y_{l-1}$ . The full recursion is derived as:

$$\alpha_{0,k} = 0 \quad \forall k \quad (17.41)$$

$$\zeta_{0,k} = \emptyset \quad \forall k \quad (17.42)$$

the score for any empty sequence is zero

$$\alpha_{l+1,k} = \max_{\hat{y}_{1:l}} w \cdot \phi_{l+1}(x, \hat{y} \circ k) \quad (17.43)$$

separate score of prefix from score of position l+1

$$= \max_{\hat{y}_{1:l}} w \cdot \left( \phi_{1:l}(x, \hat{y}) + \phi_{l+1}(x, \hat{y} \circ k) \right) \quad (17.44)$$

distributive law over dot products

$$= \max_{\hat{y}_{1:l}} \left[ w \cdot \phi_{1:l}(x, \hat{y}) + w \cdot \phi_{l+1}(x, \hat{y} \circ k) \right] \quad (17.45)$$

separate out final label from prefix, call it k'

$$= \max_{\hat{y}_{1:l-1}} \max_{k'} \left[ w \cdot \phi_{1:l}(x, \hat{y} \circ k') + w \cdot \phi_{l+1}(x, \hat{y} \circ k' \circ k) \right] \quad (17.46)$$

swap order of maxes, and last term doesn't depend on prefix

$$= \max_{k'} \left[ \left[ \max_{\hat{y}_{1:l-1}} w \cdot \phi_{1:l}(x, \hat{y} \circ k') \right] + w \cdot \phi_{l+1}(x, \langle \dots, k', k \rangle) \right] \quad (17.47)$$

apply recursive definition

$$= \max_{k'} \left[ \alpha_{l,k'} + w \cdot \phi_{l+1}(x, \langle \dots, k', k \rangle) \right] \quad (17.48)$$

**Algorithm 42** ARGMAXFORSEQUENCES( $x, w$ )

---

```

1:  $L \leftarrow \text{LEN}(x)$ 
2:  $\alpha_{l,k} \leftarrow 0, \quad \zeta_{k,l} \leftarrow 0, \quad \forall k = 1 \dots K, \quad \forall l = 0 \dots L$  // initialize variables
3: for  $l = 0 \dots L-1$  do
4:   for  $k = 1 \dots K$  do
5:      $\alpha_{l+1,k} \leftarrow \max_{k'} [\alpha_{l,k'} + w \cdot \phi_{l+1}(x, \langle \dots, k', k \rangle)]$  // recursion:
        // here,  $\phi_{l+1}(\dots, k', k, \dots)$  is the set of features associated with
        // output position  $l + 1$  and two adjacent labels  $k'$  and  $k$  at that position
6:      $\zeta_{l+1,k} \leftarrow$  the  $k'$  that achieves the maximum above // store backpointer
7:   end for
8: end for
9:  $y \leftarrow \langle 0, 0, \dots, 0 \rangle$  // initialize predicted output to L-many zeros
10:  $y_L \leftarrow \text{argmax}_k \alpha_{L,k}$  // extract highest scoring final label
11: for  $l = L-1 \dots 1$  do
12:    $y_l \leftarrow \zeta_{l, y_{l+1}}$  // traceback  $\zeta$  based on  $y_{l+1}$ 
13: end for
14: return  $y$  // return predicted output

```

---

and record a backpointer to the  $k'$  that achieves the max

$$\zeta_{l+1,k} = \underset{k'}{\text{argmax}} [\alpha_{l,k'} + w \cdot \phi_{l+1}(x, \langle \dots, k', k \rangle)] \quad (17.49)$$

At the end, we can take  $\max_k \alpha_{L,k}$  as the score of the best output sequence. To extract the final sequence, we know that the best label for the last word is  $\text{argmax}_k \alpha_{L,k}$ . Let's call this  $\hat{y}_L$ . Once we know that, the best *previous* label is  $\zeta_{L-1, \hat{y}_L}$ . We can then follow a path through  $\zeta$  back to the beginning. Putting this all together gives Algorithm 17.7.

The main benefit of Algorithm 17.7 is that it is guaranteed to exactly compute the  $\text{argmax}$  output for sequences required in the structured perceptron algorithm, *efficiently*. In particular, its runtime is  $O(LK^2)$ , which is an exponential improvement on the naive  $O(K^L)$  runtime if one were to enumerate every possible output sequence. The algorithm can be naturally extended to handle “higher order” Markov assumptions, where features depend on triples or quadruples of the output. The memoization becomes notationally cumbersome, but the algorithm remains essentially the same. In order to handle a length  $M$  Markov features, the resulting algorithm will take  $O(LK^M)$  time. In practice, it's rare that  $M > 3$  is necessary or useful.

In the case of loss-augmented search for structured SVMs (as opposed to structured perceptron), we need to include the scores coming from the loss augmentation in the dynamic program. The only thing that changes between the standard  $\text{argmax}$  solution (Algorithm 17.7, and derivation in Eq (17.48)) is that the any time an incorrect label is used, the (loss-augmented) score increases by one. Recall that in the non-loss-augmented case, we have the  $\alpha$  recursion



as:

$$\alpha_{l+1,k} = \max_{\hat{y}_{1:l}} \mathbf{w} \cdot \phi_{1:l+1}(\mathbf{x}, \hat{\mathbf{y}} \circ k) \quad (17.50)$$

$$= \max_{k'} \left[ \alpha_{l,k'} + \mathbf{w} \cdot \phi_{l+1}(\mathbf{x}, \langle \dots, k', k \rangle) \right] \quad (17.51)$$

If we define  $\tilde{\alpha}$  to be the loss-augmented score, the corresponding recursion is (differences highlighted in blue):

$$\tilde{\alpha}_{l+1,k} = \max_{\hat{y}_{1:l}} \mathbf{w} \cdot \phi_{1:l+1}(\mathbf{x}, \hat{\mathbf{y}} \circ k) + \ell_{1:l+1}^{(\text{Ham})}(\mathbf{y}, \hat{\mathbf{y}} \circ k) \quad (17.52)$$

$$= \max_{k'} \left[ \tilde{\alpha}_{l,k'} + \mathbf{w} \cdot \phi_{l+1}(\mathbf{x}, \langle \dots, k', k \rangle) \right] + \mathbf{1}[k \neq \mathbf{y}_{l+1}] \quad (17.53)$$

In other words, when computing  $\tilde{\alpha}$  in the loss-augmented case, whenever the output prediction is forced to pass through an incorrect label, the score for that cell in the dynamic program gets increased by one. The resulting algorithm is identical to Algorithm 17.7, except that Eq (17.53) is used for computing  $\alpha$ s.

## 17.8 Further Reading

TODO