

# CS 4650/7650, Lecture 11

## Discriminative and conditional models for sequence labeling

Jacob Eisenstein

September 23, 2013

### 1 Review of Hidden Markov Models

Hidden Markov Models (HMMs) are a special case of a weighted (probabilistic) finite-state transducer, which transduce from tags to words. The probability model is

$$P(\mathbf{x}, \mathbf{y}) = \prod_n P(x_n | y_n; \phi) P(y_n | y_{n-1}; \theta)$$

**Estimation** in HMMs is easy. We can use smoothed relative frequency estimation to compute the parameters as ratios of counts, e.g.

$$\phi_{\text{fish}|\text{NN}} = \frac{\text{count}(\text{fish}, \text{NN})}{\text{count}(\text{NN})}$$
$$\theta_{\text{NN}|\text{DT}} = \frac{\text{count}(\text{DT}, \text{NN})}{\text{count}(\text{DT})}$$

**Decoding** in HMMs is more interesting. We want to find  $\hat{\mathbf{y}} = \arg \max_{\mathbf{y}} P(\mathbf{y}, \mathbf{x})$ , but that involves searching  $\mathcal{O}(K^N)$  tag sequences.

The **Viterbi Algorithm** is based on a recurrence relation. Let's define

$$v_n(k) = \max_{\mathbf{y}_{1:n-1}} P(\mathbf{x}_{1:n}, y_n = k, \mathbf{y}_{1:n-1}), \quad (1)$$

which is the joint probability of the best tag sequence ending in  $y_n = k$ . Then we can apply the chain rule,

$$v_n(k) = P(x_n|y_n = k) \max_{\mathbf{y}_{1:n-1}} P(y_n = k|y_{n-1})P(\mathbf{x}_{1:n-1}, \mathbf{y}_{1:n-1}) \quad (2)$$

Now let's say that  $y_{n-1}$  takes the value  $k'$ . Then,

$$v_n(k) = P(x_n|y_n = k) \max_{k'} P(y_n = k|y_{n-1} = k')v_{n-1}(k'). \quad (3)$$

The Viterbi algorithm directly implements this recurrence relation, in the form of a lattice. I made some nice slides about this. Check them out. In the slides, we add log-probabilities rather than multiply probabilities. This is crucial for long sequences so that the probabilities don't underflow.

**The Forward Algorithm** computes the marginal probability  $P(\mathbf{x}_{1:N}) = \sum_{\mathbf{y}_{1:N}} P(\mathbf{x}_{1:N}, \mathbf{y}_{1:N})$ . It's almost the same exact thing as the Viterbi algorithm, but replace max with addition. Let's define

$$\begin{aligned} \alpha_n(k) &= \sum_{\mathbf{y}_{1:n-1}} P(\mathbf{x}_{1:n}, y_n = k, \mathbf{y}_{1:n-1}) \\ &= P(x_n|y_n = k) \sum_{\mathbf{y}_{1:n-1}} P(y_n = k|y_{n-1})P(\mathbf{x}_{1:n-1}, \mathbf{y}_{1:n-1}) \\ &= P(x_n|y_n = k) \sum_{k'} P(y_n = k|y_{n-1} = k')\alpha_{n-1}(k'). \end{aligned}$$

If we're in the log domain, we have to replace addition with log addition,  $\oplus(a, b) = \log(e^a + e^b)$ . So,

$$\alpha_n(k) = \log P(x_n|y_n = k) + \log \sum_{k'} \exp(\log P(y_n = k|y_{n-1} = k') + \alpha_{n-1}(k')) \quad (4)$$

You can use a numerical library like `scipy.misc.logsumexp` to avoid underflow in the computation of the exponent.

## 2 Tagging and features

The HMM can incorporate two sources of information:

- Word-tag probabilities, via  $P(x_n|y_n)$
- Local context, via  $P(y_n|y_{n-1})$ .

Note that there are a lot of things that we're missing here:

- **Morphology**. *Slithy toves* just kind of looks like JJ NNS, because of the apparent suffixes.
- **Capitalization**. This is especially relevant for named entity recognition, e.g., *I bought an apple* and *I bought an Apple phone*.
- **Word-specific context**. In the PTB, *this* and *these* are both tagged DT. But *this* is likely to be followed by a singular noun NN, and *these* is likely to be followed by a plural noun NNS.

In addition, we're learning to maximize the joint probability  $P(\mathbf{x}, \mathbf{y})$ , rather than to maximize the conditional probability  $P(\mathbf{y}|\mathbf{x})$ , or to minimize the error (or sum of hinge losses). Today we'll talk about methods that address these limitations.

### 3 Structured perceptron

Remember the perceptron update:

$$\hat{y} = \arg \max_y \mathbf{w}^\top \mathbf{f}(\mathbf{x}, y) \quad (5)$$

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \mathbf{f}(\mathbf{x}, y) - \mathbf{f}(\mathbf{x}, \hat{y}) \quad (6)$$

In sequence labeling, we have a **structured output**  $\mathbf{y} \in \mathcal{Y}(\mathbf{x})$ . Can we still apply the perceptron rule?

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}} \mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}) \quad (7)$$

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \mathbf{f}(\mathbf{x}, \mathbf{y}) - \mathbf{f}(\mathbf{x}, \hat{\mathbf{y}}) \quad (8)$$

This is called structured perceptron, because it learns to predict structured output  $\mathbf{y}$ . The only difference from standard perceptron is that the set of possible outputs  $\mathcal{Y}(\mathbf{x})$  is a function of the input: in the case of sequence labeling, it is the set of tag sequences with the same length as the input.

Suppose we want to apply this to POS tagging? What features might we want? Here are three:

- Word-tag features, e.g.  $\langle W : \text{slithy}, \text{JJ} \rangle$
- Adjacent tag-tag features, e.g.  $\langle T : \text{JJ}, \text{NNS} \rangle$
- Suffix-tag features, e.g.,  $\langle M : \text{-es}, \text{NNS} \rangle$

Then we can characterize the tagging DT JJ NNS of the text the slithy toves (from Jabberwocky) in terms of the following features:

$$\begin{aligned} \mathbf{f}(\text{the slithy toves, DT JJ NNS}) = & \{ \langle W : \text{the}, \text{DT} \rangle, \langle M : \emptyset, \text{DT} \rangle, \langle T : \star, \text{DT} \rangle \\ & \langle W : \text{slithy}, \text{JJ} \rangle, \langle M : \text{-thy}, \text{JJ} \rangle, \langle T : \text{DT}, \text{JJ} \rangle \\ & \langle W : \text{toves}, \text{NNS} \rangle, \langle M : \text{-es}, \text{NNS} \rangle, \langle T : \text{JJ}, \text{NNS} \rangle \\ & \langle T : \text{NNS}, \star \rangle \} \end{aligned}$$

If this is the correct tagging, then we hope to learn a set of weights  $\mathbf{w}$  such that  $\mathbf{w}^\top \mathbf{f}(\text{the slithy toves, DT JJ NNS}) > \mathbf{w}^\top \mathbf{f}(\text{the slithy toves, DT NN VBZ})$  (for example).

### 3.1 Decoding

Recall that perceptron learning requires computing  $\arg \max_{\mathbf{y}} \mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y})$ . The argmax is over exponentially many sequences. What can we do?

Suppose we restrict the feature function to the form,

$$\mathbf{f}(\mathbf{x}_{1:N}, \mathbf{y}_{1:N}) = \sum_n \mathbf{f}(\mathbf{x}_{1:N}, y_n, y_{n-1}, n) \quad (9)$$

Let's think about what this restriction says:

- We can incorporate the HMM “features” easily: the transition probabilities are covered by  $\langle y_n, y_{n-1} \rangle$ , and the emission probabilities are covered by  $\langle y_n, x_n \rangle$ . Note that the feature function needs to be parametrized by  $n$  for the emission feature.
- We can consider any part of the input  $\mathbf{x}_{1:N}$ . For example, we can have a feature that fires on  $y_n = \text{JJ}$  and  $x_{n-15} = \text{barnacle}$ . We can also do absolute indexing on  $x$ , for example, firing if  $x_1 = \text{slithy}$ .
- We can have features that consider morphological or character-level (“word shape”) features of  $x_n$ ; any function that doesn't consider tags other than  $y_n$  or  $y_{n-1}$  is okay.

- What we **cannot** do is consider arbitrary parts of  $\mathbf{y}_{1:N}$ , such as  $\langle y_1 = \text{CC} \rangle$  (e.g. penalizing sentences that start with conjunctions), or  $\langle y_{n-3} = \text{DT}, y_n = \text{CC} \rangle$ .

Given this restriction, we have

$$\mathbf{w}^\top \mathbf{f}(\mathbf{x}_{1:N}, \mathbf{y}_{1:N}) = \sum_n \mathbf{w}^\top \mathbf{f}(\mathbf{x}_{1:N}, y_n, y_{n-1}, n) \quad (10)$$

We can now apply the Viterbi algorithm to do **decoding**, with a very similar recurrence relation to what we saw before:

$$v_n(k) = \max_{\mathbf{y}_{1:n-1}} \mathbf{w}^\top \mathbf{f}(\mathbf{x}_{1:N}, \mathbf{y}_{1:n-1}, y_n = k) \quad (11)$$

$$= \max_{\mathbf{y}_{1:n}: y_n = k} \sum_m \mathbf{w}^\top \mathbf{f}(\mathbf{x}_{1:N}, y_m, y_{m-1}, m) \quad (12)$$

$$= \max_{y_{n-1}=k'} \mathbf{w}^\top \mathbf{f}(\mathbf{x}_{1:N}, y_n = k, y_{n-1} = k', n) + \sum_m^{n-1} \mathbf{w}^\top \mathbf{f}(\mathbf{x}_{1:N}, y_m = k', y_{m-1}, m) \quad (13)$$

$$= \max_{y_{n-1}=k'} \mathbf{w}^\top \mathbf{f}(\mathbf{x}_{1:N}, y_n = k, y_{n-1} = k', n) + v_{n-1}(k') \quad (14)$$

So to compute  $v_n(k)$ , we have to iterate over all  $y_{n-1} = k'$ ,

- find the features  $\mathbf{f}(\mathbf{x}_{1:N}, y_n = k, y_{n-1} = k', n)$ ,
- compute the inner product  $\mathbf{w}^\top \mathbf{f}(\mathbf{x}_{1:N}, y_n = k, y_{n-1} = k', n)$ ,
- add it to  $v_{n-1}(k')$
- take the max over all  $k'$

This only works because of the assumption that the feature function decomposes over local parts of the sequence! If we wanted a feature that considered arbitrary parts of tag sequence, there would be no way to incorporate it into the recurrence relation.

### 3.2 Example

In our example, we have

$$\begin{aligned}
 v_0(\text{DT}) &= \mathbf{w}^\top \{ \langle T : \star, \text{DT} \rangle, \langle W : \text{the}, \text{DT} \rangle, \langle M : \emptyset, \text{DT} \rangle \} \\
 v_0(\text{JJ}) &= \mathbf{w}^\top \{ \langle T : \star, \text{JJ} \rangle, \langle W : \text{the}, \text{JJ} \rangle, \langle M : \emptyset, \text{JJ} \rangle \} \\
 &\dots \\
 v_1(\text{JJ}) &= \max \quad \mathbf{w}^\top \{ \langle T : \text{DT}, \text{JJ} \rangle, \langle W : \text{slithy}, \text{JJ} \rangle, \langle M : \text{-thy}, \text{JJ} \rangle \} + v_0(\text{DT}), \\
 &\quad \mathbf{w}^\top \{ \langle T : \text{JJ}, \text{JJ} \rangle, \langle W : \text{slithy}, \text{JJ} \rangle, \langle M : \text{-thy}, \text{JJ} \rangle \} + v_0(\text{JJ}), \\
 &\dots
 \end{aligned}$$

### 3.3 Learning

If we define  $\mathbf{f}(\mathbf{x}_{1:N}, y_n, y_{n-1}, n) = \{ \langle W : x_n, y_n \rangle, \langle T : y_{n-1}, y_n \rangle \}$ , then our model is identical to the HMM. If we set the weights of these features to the log of their maximum-likelihood estimates,

$$\begin{aligned}
 w_{\langle W : x_n, y_n \rangle} &= \log \text{count}(x_n, y_n) - \log \text{count}(y_n) \\
 w_{\langle W : y_{n-1}, y_n \rangle} &= \log \text{count}(y_{n-1}, y_n) - \log \text{count}(y_{n-1}),
 \end{aligned}$$

then we exactly recover the HMM.

But to use more overlapping features and to get the advantages of error-driven learning, we're going to do perceptron updates. It's exactly the same as the non-structured perceptron:

- $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \mathbf{f}(\mathbf{x}, \mathbf{y}) - \mathbf{f}(\mathbf{x}, \hat{\mathbf{y}})$  is the standard update, using Viterbi to find  $\hat{\mathbf{y}}$ .
- As before, we can use MIRA to do large-margin training.

## 4 Conditional random fields

Structured perceptron works well in practice, and you will implement in your project 2, where it should work much much better than the Hidden Markov Model.

- But sometimes we need probabilities, and SP doesn't give us that.

- The Conditional Random Field (CRF) is a probabilistic conditional model for sequence labeling.
- Just as structured perceptron is built on the perceptron classifier, conditional random fields are built on the logistic regression classifier.

$$P(y|x) = \frac{e^{\mathbf{w}^\top \mathbf{f}(y,x)}}{\sum_{y' \in \mathcal{Y}} e^{\mathbf{w}^\top \mathbf{f}(y',x)}}$$

We can again move to structured prediction, assuming the same restriction on the scoring function:

$$P(\mathbf{y}|\mathbf{x}) = \frac{e^{\mathbf{w}^\top \mathbf{f}(\mathbf{y},\mathbf{x})}}{\sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{x})} e^{\mathbf{w}^\top \mathbf{f}(\mathbf{y}',\mathbf{x})}}$$

This is called a conditional random field because it models the sequence labeling task as a Markov random field, and estimates the probability of a set of variables **conditioned** on the others (as opposed to jointly, which the HMM does).

## 4.1 Decoding in CRFs

Decoding in the CRF does not depend on the denominator.

$$\begin{aligned} \hat{\mathbf{y}} &= \arg \max_{\mathbf{y}} P(\mathbf{y}|\mathbf{x}) \\ &= \arg \max_{\mathbf{y}} \log P(\mathbf{y}|\mathbf{x}) \\ &= \arg \max_{\mathbf{y}} \mathbf{w}^\top \mathbf{f}(\mathbf{y}, \mathbf{x}) - \log \sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{x})} e^{\mathbf{w}^\top \mathbf{f}(\mathbf{y}', \mathbf{x})} \\ &= \arg \max_{\mathbf{y}} \mathbf{w}^\top \mathbf{f}(\mathbf{y}, \mathbf{x}) \\ &= \arg \max_{\mathbf{y}} \sum_n^N \mathbf{w}^\top \mathbf{f}_n(y_n, y_{n-1}, \mathbf{x}_{1:N}, n) \end{aligned}$$

So we can just apply Viterbi directly, as defined in the structured perceptron. At each  $n$ , we have,

$$\begin{aligned} v_n(k) &= \max_{k'} v_{n-1}(k') + \mathbf{w}^\top \mathbf{f}(y_n = k, y_{n-1} = k', \mathbf{x}, n) \\ &= \bigoplus_{k'} v_{n-1}(k') \otimes \mathbf{w}^\top \mathbf{f}(y_n = k, y_{n-1} = k', \mathbf{x}, n) \end{aligned}$$

## 4.2 Learning in CRFs

Learning is a little more complicated. As with logistic regression, we need to learn weights to maximize the conditional log probability,

$$\begin{aligned}\ell &= \sum_i^{\text{\#instances}} \log P(\mathbf{y}_i | \mathbf{x}_i), \\ &= \sum_i \mathbf{w}^\top \mathbf{f}(\mathbf{x}_i, \mathbf{y}_i) - \log \sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{x}_i)} \exp(\mathbf{w}^\top \mathbf{f}(\mathbf{x}_i, \mathbf{y}'))\end{aligned}$$

As in logistic regression, the derivative is a difference between observed and expected counts:

$$\begin{aligned}\frac{d\ell}{dw_j} &= \sum_i \text{count}(\mathbf{x}_i, \mathbf{y}_i)_j - E_{\mathbf{y}|\mathbf{x}_i; \mathbf{w}}[\text{count}(\mathbf{x}_i, \mathbf{y})_j] \\ \text{count}(\mathbf{x}_i, \mathbf{y}_i)_j &= \sum_n^N f_{n,j}(\mathbf{x}_i, y_{i,n}, y_{i,n-1}, n)\end{aligned}$$

For example:

- If feature  $j$  is  $\langle CC, DT \rangle$ , then  $c_j(\mathbf{x}_n, \mathbf{y}_n)$  is the count of times DT follows CC in the sequence  $\mathbf{y}_n$ .
- If feature  $j$  is  $\langle M : -thy, JJ \rangle$ , then  $\text{count}(\mathbf{x}_n, \mathbf{y}_n)_j$  is the count of words ending in *-thy* in  $\mathbf{x}_n$  that are tagged JJ.

The expected feature counts are more complex.

- $E_{\mathbf{y}|\mathbf{x}; \mathbf{w}}[\text{count}(\mathbf{x}_i, \mathbf{y})_j] = \sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{x}_i)} P(\mathbf{y} | \mathbf{x}_i; \mathbf{w}) f_j(\mathbf{x}, \mathbf{y})$
- This looks bad: we have to sum over an exponential number of labelings again.



- But remember that the feature function decomposes  $f_j(\mathbf{x}, \mathbf{y}) = \sum_n f_j(\mathbf{x}, y_n, y_{n-1}, n)$ .

$$\begin{aligned}
E_{\mathbf{y}|\mathbf{x};\mathbf{w}}[\text{count}(\mathbf{x}, \mathbf{y})_j] &= \sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} P(\mathbf{y}|\mathbf{x}; \mathbf{w}) f_j(\mathbf{x}, \mathbf{y}) \\
&= \sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} P(\mathbf{y}|\mathbf{x}; \mathbf{w}) \sum_n^N f_j(\mathbf{x}, y_n, y_{n-1}, n) \\
&= \sum_n^N \sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} P(\mathbf{y}|\mathbf{x}; \mathbf{w}) f_j(\mathbf{x}, y_n, y_{n-1}, n) \\
&= \sum_n^N \sum_{j,k \in \mathcal{Y}} \sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{x}): y_{n-1}=j, y_n=k} P(\mathbf{y}|\mathbf{x}; \mathbf{w}) f_j(\mathbf{x}, y_n, y_{n-1}, n) \\
&= \sum_n^N \sum_{j,k \in \mathcal{Y}} f_j(\mathbf{x}, y_n, y_{n-1}, n) \sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{x}): y_{n-1}=j, y_n=k} P(\mathbf{y}|\mathbf{x}; \mathbf{w}) \\
&= \sum_n^N \sum_{j,k \in \mathcal{Y}} f_j(\mathbf{x}, y_n, y_{n-1}, n) P(y_{n-1} = j, y_n = k | \mathbf{x}; \mathbf{w})
\end{aligned}$$

- The expected feature counts can be computed efficiently if we know the **marginal** probabilities  $P(y_n, y_{n-1} | \mathbf{x}; \mathbf{w})$ .
- This is the probability of traversing the edge  $y_{n-1} \rightarrow y_n$ , conditioned on the entire observation  $\mathbf{x}_{1:N}$ . [\[Draw this in trellis\]](#)
- To compute this marginal probability, we will apply the forward-backward algorithm.

## 5 The forward-backward algorithm

Here we require the marginal probability, e.g.  $P(y_n = \text{NNP}, y_{n-1} = \text{DET} | x_{1:N})$ .

We can rewrite the marginal probability of a single tag transition  $P(y_n, y_{n-1} | x_{1:N})$  as a ratio:

$$P(y_n, y_{n-1} | x_{1:N}) = \frac{P(y_n, y_{n-1}, x_{1:N})}{P(x_{1:N})} \quad (15)$$

Recall the Forward Algorithm, which defines the forward probabilities  $\alpha_n(k) \triangleq P(y_n, x_{1:n})$ . This allows us to compute the denominator  $P(x_{1:N}) = \sum_k \alpha_N(k)$ .

## 5.1 Forward algorithm

First we need to show how to use the forward algorithm in a CRF. We're going to switch to a semiring in which we work with probabilities rather than log-probabilities, so that we can use simple addition rather than log-addition. But in practice we'd work in a log-probability semiring, where  $a \oplus b = \log(e^a + e^b)$ . For a practical implementation, see `numpy.logaddexp` and `scipy.misc.logsumexp`.

$$\alpha_n(k) \triangleq P(y_n = k, x_{1:n}) = \sum_{\mathbf{y}_{1:n}: y_n = k} P(\mathbf{y}_{1:n}, \mathbf{x}_{1:n}) \quad (16)$$

$$= \sum_{\mathbf{y}_{1:n}: y_n = k} \exp\left(\sum_{m=1}^n \mathbf{w}^\top \mathbf{f}(y_m, y_{m-1}, \mathbf{x}, m)\right) \quad (17)$$

$$= \sum_{\mathbf{y}_{1:n}} \exp(\mathbf{w}^\top \mathbf{f}(y_n = k, y_{n-1}, \mathbf{x}, n)) \exp\left(\sum_{m=1}^{n-1} \mathbf{w}^\top \mathbf{f}(y_m, y_{m-1}, \mathbf{x}, m)\right) \quad (18)$$

$$= \sum_{y_{n-1}=k'} \exp(\mathbf{w}^\top \mathbf{f}(y_n = k, y_{n-1} = k', \mathbf{x}, n)) \quad (19)$$

$$\times \sum_{\mathbf{y}_{1:n-1}: y_{n-1} = k'} \exp\left(\sum_{m=1}^{n-1} \mathbf{w}^\top \mathbf{f}(y_m, y_{m-1}, \mathbf{x}, m)\right) \quad (20)$$

$$= \sum_{y_{n-1}=k'} \exp(\mathbf{w}^\top \mathbf{f}(y_n = k, y_{n-1} = k', \mathbf{x}, n)) \alpha_{n-1}(k') \quad (21)$$

$$P(\mathbf{x}_{1:N}) = \sum_k \alpha_N(y_N = k) \quad (22)$$

We can see an example of this in the slides. In this example, we have only emission and transition features. The feature weights are shown in the

slide. So to compute the forward variable for *can* tagged as N, we have,

$$\alpha_2(N) = (-4 \otimes -3 \otimes -3) \oplus (-5 \otimes -1 \otimes -3) \quad (23)$$

$$= (-10) \oplus (-9) \quad (24)$$

$$= \log(e^{-10} + e^{-9}) \quad (25)$$

$$\approx -8.7 \quad (26)$$

Now we return to the numerator of  $P(y_n, y_{n-1} | x_{1:N}) = \frac{P(y_n, y_{n-1}, x_{1:N})}{P(x_{1:N})}$ . We can factor,

$$\begin{aligned} P(y_n, y_{n-1}, x_{1:N}) &= P(\mathbf{x}_{n+1:N} | y_n, y_{n-1}, x_{1:N}) P(y_n, y_{n-1}, x_{1:N}) \\ &= P(\mathbf{x}_{n+1:N} | y_n) P(y_n, x_n | y_{n-1}) P(y_{n-1}, x_{1:n-1}) \\ &\propto P(\mathbf{x}_{n+1:N} | y_n) \exp(\mathbf{w}^\top \mathbf{f}(y_n, y_{n-1}, \mathbf{x}, n)) \alpha_{n-1}(y_{n-1}), \end{aligned}$$

where we have plugged in the forward variable  $\alpha_n(y_{n-1}) = P(y_{n-1}, \mathbf{x}_{1:n-1})$ . Now we define a **backward variable**,  $b_n(y_n) = P(\mathbf{x}_{n+1:N} | y_n)$ , yielding

$$P(y_n, y_{n-1} | x_{1:N}) \propto b_n(y_n) \exp(\mathbf{w}^\top \mathbf{f}(y_n, y_{n-1}, \mathbf{x}, n)) \alpha_{n-1}(y_{n-1}).$$

If we want to compute the marginal for a single node, we have

$$P(y_n | \mathbf{x}_{1:N}) \propto b_n(y_n) \alpha_n(y_n) = P(\mathbf{x}_{n+1:N} | y_n) P(y_n, x_{1:n}) = P(y_n, \mathbf{x}_{1:N}) \quad (27)$$

This is useful for applying expectation-maximization to unsupervised learning of HMMs.

## 5.2 The backwards algorithm

We will need another recursive algorithm to compute the backwards variables. Fortunately, the backwards probability also decomposes into a recurrence relation,

$$\begin{aligned} b_n(y_n) &\triangleq P(x_{n+1:N} | y_n) \\ &= \sum_{y_{n+1}} P(x_{n+1:N}, y_{n+1} | y_n) \\ &= \sum_{y_{n+1}} P(x_{n+2:N} | y_{n+1}, x_{n+1}) P(x_{n+1}, y_{n+1} | y_n) \\ &\propto \sum_{y_{n+1}} b_{n+1}(y_{n+1}) \exp(\mathbf{w}^\top \mathbf{f}(y_{n+1}, y_n, \mathbf{x}_{1:N}, n)) \\ b_N(y_N) &= P(\text{END} | y_N) \end{aligned}$$

The steps are as follows:

- introduce  $y_{n+1}$  and sum over all possible values
- apply the chain rule to extract  $y_{n+1}$  and  $x_{n+1}$
- use the definition of the backwards probability and the model assumptions

The backwards algorithm moves from right to left, starting with  $B_N(y_N) = P(\text{END}|y_N)$ , the score for transitioning to the end state given  $y_N$ .

We can see an example in the slides.

### 5.3 Learning in CRFs: wrapup

The overall procedure looks just like logistic regression:

- Use forward backward to compute expected feature counts under  $P(\mathbf{y}|\mathbf{x}; \mathbf{w})$
- Compute gradients as difference between feature counts and expected counts
- Update  $\mathbf{w}$  (using quasinewton optimization)
- Iterate