# Yet Another Introduction to Hidden Markov Models

**Anoop Sarkar**
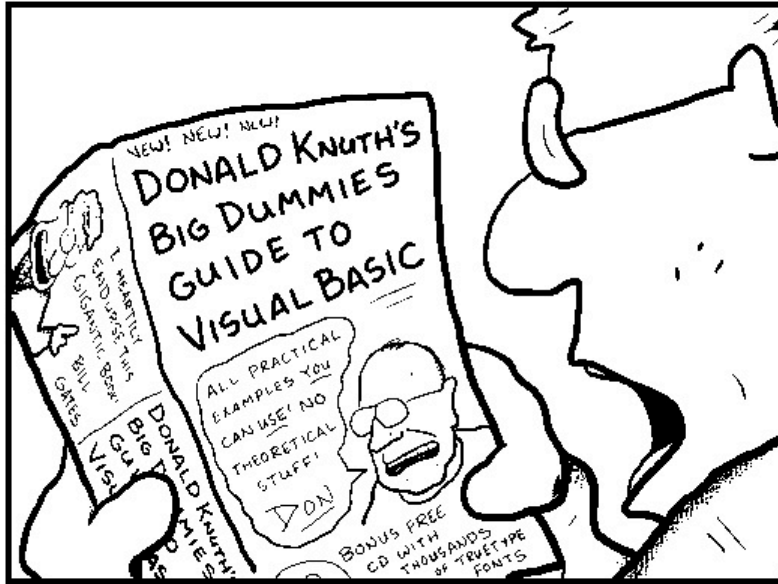School of Computing Science, Simon Fraser University
Burnaby, BC V5A 1S6, Canada
anoop@cs.sfu.ca

## DOCTOR FUN

10 Feb 2000

Don Knuth finally sells out.

## Abstract

Yet another introduction to Baum-Welch re-estimation algorithm for learning parameters of a Hidden Markov Model (HMM). The presentation is from an NLP perspective and the focus is on deriving efficient pseudo-code for a fast implementation of the Viterbi, forward and backward algorithms that are used in Baum-Welch re-estimation. We consider only discrete states and discrete observations. In fact, to make the presentation NLP friendly we assume the task will be to train from word sequences with tagging sequences as hidden data, where in the decoding step each word will be eventually tagged with its appropriate tag. The tags are discrete, but can refer to part of speech tags, chunking tags, tags that indicate word boundaries, and other NLP applications that can be represented as a sequence learning problem.

## 1 Notation

We refer to states or tags as $t \in \{1, \ldots, \mathtt{N}\}$ and observations or words as $w \in \{1, \ldots, \mathtt{M}\}$. In place of the tag index or word index we sometimes use an actual tag $t = \mathtt{DT}$ or word $w = \mathtt{the}$ in some examples. A special word $w = 0$ and tag $t = 0$ is used to indicate the begin and end of a sequence. The sequence of states and associated observations, or tags and associated words is $w_0, t_0, \ldots, w_{T-1}, t_{T-1}$, a sequence of length $T$. We first consider a single sequence, later generalizing to multiple sequences of varying length.

## 2 Defining a Hidden Markov Model

A Hidden Markov Model (HMM) $\theta$ is the triple $\langle p_{\mathsf{s}}, p_{\mathsf{tt}}, p_{\mathsf{tw}} \rangle$, where

1. $p_{\mathsf{s}}(t_0)$ is the probability that we start with some tag $t$ as $t_0$,

2. $p_{\mathsf{tt}}(t_i \mid t_{i-1})$ is the transition probability from $t_{i-1}$ to $t_i$, and

3. $p_{\mathsf{tw}}(w_i \mid t_i)$ is the probability of generating $w_i$ from $t_i$.

We set $p_{\mathsf{s}}(0) = 1.0$ and $p_{\mathsf{tw}}(0 \mid 0) = 1.0$. As a result we can now ignore $p_{\mathsf{s}}$ and our HMM is represented by the tuple $\langle p_{\mathsf{tt}}, p_{\mathsf{tw}} \rangle$.

## 3 The Forward Algorithm

Forward step ($\alpha$):

$$
\begin{aligned}
\Pr(w_0 = 0 \mid t_0 = 0) &= 1.0 \\
\Pr(w_0, \ldots, w_i, t_i) &= \sum_{t_{i-1}} \Pr(w_0, \ldots, w_{i-1}, t_{i-1}) \cdot p_{\mathsf{tt}}(t_i \mid t_{i-1}) \cdot p_{\mathsf{tw}}(w_i \mid t_i) \\
\Pr(w_0, \ldots, w_{T-1}) &= \sum_{t_{T-1}} \Pr(w_0, \ldots, w_{T-1}, t_{T-1}) \\
&= \Pr(w_0, \ldots, w_{T-1}, t_{T-1} = 0)
\end{aligned}
\tag{1}
$$

$$
\begin{aligned}
\alpha(0, 0) &= 1.0 \\
\alpha(t_i, i) &= \sum_{t_{i-1}} \alpha(t_{i-1}, i-1) \cdot p_{\mathsf{tt}}(t_i \mid t_{i-1}) \cdot p_{\mathsf{tw}}(w_i \mid t_i), \quad 1 \le i \le T-1 \\
\Pr(w_0, \ldots, w_{T-1}) &= \sum_{t} \alpha(t, T-1) \\
&= \alpha(0, T-1)
\end{aligned}
\tag{2}
$$

**Algorithm 1** $\alpha$: implements the forward algorithm

---

**Require:** $W = w_0, \ldots, w_{T-1}$
**Require:** *tag-dict*: $\{1, \ldots, \texttt{M}\} \to \mathcal{P}(\{1, \ldots, \texttt{N}\})$
  1: $\alpha(0, 0) := 1.0$
  2: **for** $(1 \le i \le T - 1)$ **do**
  3:     **for** $(t_i \in \textit{tag-dict}(w_i))$ **do**
  4:         **for** $(t_{i-1} \in \textit{tag-dict}(w_{i-1}))$ **do**
  5:             $\alpha(t_i, i) := \alpha(t_i, i) + \alpha(t_{i-1}, i - 1) \cdot p_{\mathsf{tt}}(t_i \mid t_{i-1}) \cdot p_{\mathsf{tw}}(w_i \mid t_i)$
  6:         **end for**
  7:     **end for**
  8: **end for**
  9: $s := \alpha(0, T - 1)$
10: **return** $s$

---

**Algorithm 2** $\alpha$: implements the forward algorithm, 1st optimization

---

**Require:** $W = w_0, \ldots, w_{T-1}$
**Require:** *tag-dict*: $\{1, \ldots, \texttt{M}\} \to \mathcal{P}(\{1, \ldots, \texttt{N}\})$
  1: $\alpha(0, 0) := 1.0$
  2: **for** $(1 \le i \le T - 1)$ **do**
  3:     **for** $(t_i \in \textit{tag-dict}(w_i))$ **do**
  4:         **for** $(t_{i-1} \in \textit{tag-dict}(w_{i-1}))$ **do**
  5:             $\alpha(t_i, i) := \alpha(t_i, i) + \alpha(t_{i-1}, i - 1) \cdot p_{\mathsf{tt}}(t_i \mid t_{i-1})$
  6:         **end for**
  7:         $\alpha(t_i, i) := \alpha(t_i, i) \cdot p_{\mathsf{tw}}(w_i \mid t_i)$
  8:     **end for**
  9: **end for**
10: $s := \alpha(0, T - 1)$
11: **return** $s$

---

## 4 The Backward Algorithm

Backward step ($\beta$):

$$
\begin{aligned}
\Pr(w_{T-1} = 0 \mid t_{T-2}) &= 1.0 \quad \text{for all } t_{T-2} \\
\Pr(w_i, \ldots, w_{T-1} \mid t_{i-1}) &= \sum_{t_i} \Pr(w_{i+1}, \ldots, w_{T-1} \mid t_i) \cdot p_{\mathsf{tt}}(t_i \mid t_{i-1}) \cdot p_{\mathsf{tw}}(w_i \mid t_i) \\
\Pr(w_0, \ldots, w_{T-1}) &= \sum_{t_0} \Pr(w_1, \ldots, w_{T-1} \mid t_0) \\
&= \Pr(w_1, \ldots, w_{T-1} \mid 0)
\end{aligned}
\tag{3}
$$

$$
\begin{aligned}
\beta(0, T - 1) &= 1.0 \\
\beta(t_{i-1}, i - 1) &= \sum_{t_i} \beta(t_i, i) \cdot p_{\mathsf{tt}}(t_i \mid t_{i-1}) \cdot p_{\mathsf{tw}}(w_i \mid t_i), \quad T - 1 \ge i > 0 \\
\Pr(w_0, \ldots, w_{T-1}) &= \sum_t \beta(t, 0) \\
&= \beta(0, 0)
\end{aligned}
\tag{4}
$$

**Algorithm 3** $\beta$: implements the backward algorithm

---

**Require:** $W = w_0, \ldots, w_{T-1}$
**Require:** *tag-dict*: $\{1, \ldots, \texttt{M}\} \to \mathcal{P}(\{1, \ldots, \texttt{N}\})$
 1: $\beta(0, T-1) := 1.0$
 2: **for** $(T-1 \geq i \geq 0)$ **do**
 3:    **for** $(t_i \in \textit{tag-dict}(w_i))$ **do**
 4:       **for** $(t_{i-1} \in \textit{tag-dict}(w_{i-1}))$ **do**
 5:          $\beta(t_{i-1}, i-1) := \beta(t_{i-1}, i-1) + p_{\mathsf{tt}}(t_i \mid t_{i-1}) \cdot p_{\mathsf{tw}}(w_i \mid t_i) \cdot \beta(t_i, i)$
 6:       **end for**
 7:    **end for**
 8: **end for**

---

## 5   Baum-Welch Re-estimation

Computing $\gamma$:

$$
\begin{aligned}
\gamma(i, t_i) &= \Pr(t_i \mid w_0, \ldots, w_{T-1}) \quad \text{for } t_i \in \{1, \ldots, \texttt{N}\} \\
&= \frac{\Pr(w_0, \ldots, w_i, t_i) \cdot \Pr(w_{i+1}, \ldots, w_{T-1} \mid t_i)}{\sum_{s_i} \Pr(w_0, \ldots, w_i, s_i) \cdot \Pr(w_{i+1}, \ldots, w_{T-1} \mid s_i)} \\
&= \frac{\Pr(w_0, \ldots, w_i, t_i) \cdot \Pr(w_{i+1}, \ldots, w_{T-1} \mid t_i)}{\Pr(w_0, \ldots, w_{T-1})} \\
&= \frac{\alpha(t_i, i) \cdot \beta(t_i, i)}{\alpha(0, T-1)}
\end{aligned}
\tag{5}
$$

Computing $\xi$:

$$
\begin{aligned}
\xi(i, t_i, t_{i+1}) &= \Pr(t_i, t_{i+1} \mid w_0, \ldots, w_{T-1}) \quad \text{for } t_i, t_{i+1} \in \{1, \ldots, \texttt{N}\} \\
&= \frac{\Pr(w_0, \ldots, w_i, t_i) \cdot p_{\mathsf{tt}}(t_{i+1} \mid t_i) \cdot p_{\mathsf{tw}}(w_{i+1} \mid t_{i+1}) \cdot \Pr(w_{i+2}, \ldots, w_{T-1} \mid t_{i+1})}{\sum_{s_i, s_{i+1}} \Pr(w_0, \ldots, w_i, s_i) \cdot p_{\mathsf{tt}}(s_{i+1} \mid s_i) \cdot p_{\mathsf{tw}}(w_{i+1} \mid s_{i+1}) \cdot \Pr(w_{i+2}, \ldots, w_{T-1} \mid s_{i+1})} \\
&= \frac{\alpha(t_i, i) \cdot p_{\mathsf{tt}}(t_{i+1} \mid t_i) \cdot p_{\mathsf{tw}}(w_{i+1} \mid t_{i+1}) \cdot \beta(t_{i+1}, i+1)}{\sum_{s_i, s_{i+1}} \alpha(s_i, i) \cdot p_{\mathsf{tt}}(s_{i+1} \mid s_i) \cdot p_{\mathsf{tw}}(w_{i+1} \mid s_{i+1}) \cdot \beta(s_{i+1}, i+1)} \\
&= \frac{\alpha(t_i, i) \cdot p_{\mathsf{tt}}(t_{i+1} \mid t_i) \cdot p_{\mathsf{tw}}(w_{i+1} \mid t_{i+1}) \cdot \beta(t_{i+1}, i+1)}{\Pr(w_0, \ldots, w_{T-1})} \\
&= \frac{\alpha(t_i, i) \cdot p_{\mathsf{tt}}(t_{i+1} \mid t_i) \cdot p_{\mathsf{tw}}(w_{i+1} \mid t_{i+1}) \cdot \beta(t_{i+1}, i+1)}{\alpha(0, T-1)}
\end{aligned}
\tag{6}
$$

To save space in computing $\xi$, we compute $\xi'$:

$$
\begin{aligned}
\xi'(t_i, t_{i+1}) &= \sum_{i=0}^{T-2} \frac{\alpha(t_i, i) \cdot p_{\mathsf{tt}}(t_{i+1} \mid t_i) \cdot p_{\mathsf{tw}}(w_{i+1} \mid t_{i+1}) \cdot \beta(t_{i+1}, i+1)}{\alpha(0, T-1)} \\
&\quad \text{for all } t_i, t_{i+1} \in \{1, \ldots, \texttt{N}\}
\end{aligned}
\tag{7}
$$

and $\delta'$:

$$
\delta'(t, w) = \sum_{i=0}^{T-1} \gamma(i, t_i = t) \cdot I(w_i = w)
$$

$$\text{where } I(\texttt{true}) = 1 \text{ and } I(\texttt{false}) = 0.$$

$$\text{for all } t \in \{1, \ldots, \texttt{N}\}, w \in \{1, \ldots, \texttt{M}\} \tag{8}$$

The new estimates for $p_{\text{tt}}$ and $p_{\text{tw}}$ are:

$$p_{\text{tt}}(t_{i+1} \mid t_i) = \frac{\xi'(t_i, t_{i+1})}{\sum_{s_i} \xi'(t_i, s_i)} \tag{9}$$

$$p_{\text{tw}}(w \mid t_i) = \frac{\delta'(t_i, w)}{\sum_{w'} \delta'(t_i, w')} \tag{10}$$

---

**Algorithm 4** Forward-Backward:
1 iteration of iterative algorithm to get new values for $p_{\text{tt}}$ and $p_{\text{tw}}$,
computes $\beta$ at the same time

---

**Require:** $W = w_0, \ldots, w_{T-1}$
**Require:** *tag-dict*: $\{1, \ldots, \texttt{M}\} \rightarrow \mathcal{P}(\{1, \ldots, \texttt{N}\})$
1:   $s := \alpha(W, \textit{tag-dict})$
2:   $\beta(0, T-1) := 1.0$
3:   **for** $(T - 1 \geq i \geq 0)$ **do**
4:      **for** $(t_i \in \textit{tag-dict}(w_i))$ **do**
5:        $\delta'(t_i, w_i) := \delta'(t_i, w_i) + \frac{1}{s} \cdot \alpha(i, t_i) \cdot \beta(i, t_i)$
6:        $\delta''(t_i) := \delta''(t_i) + \delta'(t_i, w_i)$
7:        **for** $(t_{i-1} \in \textit{tag-dict}(w_{i-1}))$ **do**
8:          $p := p_{\text{tt}}(t_i \mid t_{i-1}) \cdot p_{\text{tw}}(w_i \mid t_i)$
9:          $\beta(t_{i-1}, i-1) := \beta(t_{i-1}, i-1) + p \cdot \beta(t_i, i)$
10:         $\xi'(t_{i-1}, t_i) := \xi'(t_{i-1}, t_i) + \frac{1}{s} \cdot \alpha(t_{i-1}, i-1) \cdot p \cdot \beta(t_i, i)$
11:         $\xi''(t_{i-1}) := \xi''(t_{i-1}) + \xi'(t_{i-1}, t_i)$
12:        **end for**
13:      **end for**
14: **end for**
15: **for** $(t_i \in \{1, \ldots, \texttt{N}\})$ **do**
16:      **for** $(w \in \{1, \ldots, \texttt{M}\})$ **do**
17:        $p_{\text{tw}}(w \mid t_i) := \frac{\delta'(t_i, w)}{\delta''(t_i)}$
18:      **end for**
19:      **for** $(t_{i+1} \in \{1, \ldots, \texttt{N}\})$ **do**
20:        $p_{\text{tt}}(t_{i+1} \mid t_i) := \frac{\xi'(t_i, t_{i+1})}{\xi''(t_i)}$
21:      **end for**
22: **end for**

---

# 6   Conclusion

Did you understand it?

# References

Jason Eisner. 2004. *Tagging with a Hidden Markov Model. Lecture Notes from Johns Hopkins course 600.465 — Introduction to NLP*

John Langford. 2000. *Optimizing Hidden Markov Model Learning.* unpublished manuscript.

**Algorithm 5** Forward-Backward: 1st optimization

**Require:** $W = w_0, \ldots, w_{T-1}$
**Require:** *tag-dict*: $\{1, \ldots, \texttt{M}\} \to \mathcal{P}(\{1, \ldots, \texttt{N}\})$
 1: $s := \alpha(W, \textit{tag-dict})$
 2: $\beta(0, T - 1) := 1.0$
 3: **for** $(T - 1 \geq i \geq 0)$ **do**
 4:     **for** $(t_i \in \textit{tag-dict}(w_i))$ **do**
 5:         $\delta'(t_i, w_i) := \delta'(t_i, w_i) + \frac{1}{s} \cdot \alpha(i, t_i) \cdot \beta(i, t_i)$
 6:         $\delta''(t_i) := \delta''(t_i) + \delta'(t_i, w_i)$
 7:         $m := p_{\mathsf{tw}}(w_i \mid t_i) \cdot \beta(t_i, i)$
 8:         **for** $(t_{i-1} \in \textit{tag-dict}(w_{i-1}))$ **do**
 9:            $n := p_{\mathsf{tt}}(t_i \mid t_{i-1}) \cdot m$
10:            $\beta(t_{i-1}, i - 1) := \beta(t_{i-1}, i - 1) + n$
11:            $\xi'(t_{i-1}, t_i) := \xi'(t_{i-1}, t_i) + \frac{1}{s} \cdot \alpha(t_{i-1}, i - 1) \cdot n$
12:            $\xi''(t_{i-1}) := \xi''(t_{i-1}) + \xi'(t_{i-1}, t_i)$
13:         **end for**
14:     **end for**
15: **end for**
16: **for** $(t_i \in \{1, \ldots, \texttt{N}\})$ **do**
17:     **for** $(w \in \{1, \ldots, \texttt{M}\})$ **do**
18:         $p_{\mathsf{tw}}(w \mid t_i) := \frac{\delta'(t_i, w)}{\delta''(t_i)}$
19:     **end for**
20:     **for** $(t_{i+1} \in \{1, \ldots, \texttt{N}\})$ **do**
21:         $p_{\mathsf{tt}}(t_{i+1} \mid t_i) := \frac{\xi'(t_i, t_{i+1})}{\xi''(t_i)}$
22:     **end for**
23: **end for**

---

**Algorithm 6** Forward-Backward: multiple training sentences

---

**Require:** $S = S_0, \ldots, S_{K-1}$, a list of training sentences
**Require:** *tag-dict*: $\{1, \ldots, \mathtt{M}\} \to \mathcal{P}(\{1, \ldots, \mathtt{N}\})$

1: **for** $(0 \leq k \leq K-1)$ **do**
2:  $S_k = w_0, \ldots, w_{T-1}$
3:  $s := \alpha(S_k, \textit{tag-dict})$
4:  $L := L \cdot s$
5:  $\beta(0, T-1) := 1.0$
6:  **for** $(T-1 \geq i \geq 0)$ **do**
7:   **for** $(t_i \in \textit{tag-dict}(w_i))$ **do**
8:    $\delta'(t_i, w_i) := \delta'(t_i, w_i) + \frac{1}{s} \cdot \alpha(i, t_i) \cdot \beta(i, t_i)$
9:    $\delta''(t_i) := \delta''(t_i) + \delta'(t_i, w_i)$
10:    $m := p_{\mathsf{tw}}(w_i \mid t_i) \cdot \beta(t_i, i)$
11:    **for** $(t_{i-1} \in \textit{tag-dict}(w_{i-1}))$ **do**
12:     $n := p_{\mathsf{tt}}(t_i \mid t_{i-1}) \cdot m$
13:     $\beta(t_{i-1}, i-1) := \beta(t_{i-1}, i-1) + n$
14:     $\xi'(t_{i-1}, t_i) := \xi'(t_{i-1}, t_i) + \frac{1}{s} \cdot \alpha(t_{i-1}, i-1) \cdot n$
15:     $\xi''(t_{i-1}) := \xi''(t_{i-1}) + \xi'(t_{i-1}, t_i)$
16:    **end for**
17:   **end for**
18:  **end for**
19: **end for**
20: **for** $(t_i \in \{1, \ldots, \mathtt{N}\})$ **do**
21:  **for** $(w \in \{1, \ldots, \mathtt{M}\})$ **do**
22:   $p_{\mathsf{tw}}(w \mid t_i) := \frac{\delta'(t_i, w)}{\delta''(t_i)}$
23:  **end for**
24:  **for** $(t_{i+1} \in \{1, \ldots, \mathtt{N}\})$ **do**
25:   $p_{\mathsf{tt}}(t_{i+1} \mid t_i) := \frac{\xi'(t_i, t_{i+1})}{\xi''(t_i)}$
26:  **end for**
27: **end for**
28: **return** $L$, the likelihood of the training set $S$

---

---

**Algorithm 7** Baum-Welch Re-estimation: using only unlabelled data

---

**Require:** $S = S_0, \ldots, S_{K-1}$, a list of unlabelled training sentences
**Require:** *threshold*: small positive number, e.g. $10^{-5}$, change in likelihood for stopping condition

1: **for** $(w \in S)$ **do**
2:  $\textit{tag-dict}(w) := \{1, \ldots, \mathtt{N}\}$
3: **end for**
4: $L := 0$
5: **repeat**
6:  $L' := L$
7:  $L := \textit{Forward-Backward}(S, \textit{tag-dict})$
8:  **assert** $(L - L') > 0$
9: **until** $(L - L') > \textit{threshold}$
10: **return** $p_{\mathsf{tt}}, p_{\mathsf{tw}}$

---