

Algorithms for Hidden Markov Models continued...

**A. Overview**

- (1) a. Forward probabilities
- b. Backward probabilities
- c. Training a simple Markov chain
- d. Baum-Welch

**B. Forward probabilities**

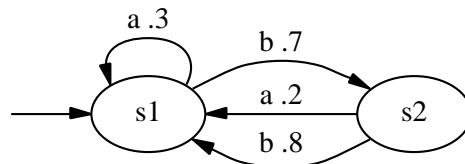
- (2)  $\alpha_i(t)$  the probability of producing  $w_{1,t-1}$  and ending up in  $s^i$ .
- (3)  $\alpha_i(t) \stackrel{\text{def}}{=} p(w_{1,t-1}, S_t = s^i)$ , where  $t > 1$
- (4)  $\alpha_i(1) = \pi_i$
- (5)  $\alpha_j(t + 1) = \sum_{i=1}^{\sigma} \alpha_i(t) p(s^i \xrightarrow{w_t} s^j)$

**C. Backward probabilities**

- (6)  $\beta_i(t)$  is the overall probability of producing  $w_{t,n}$  where the HMM is in state  $s^i$  at time  $t$ .
- (7)  $\beta_i(n + 1) = p(\epsilon | S_{n+1} = s^i) = 1$
- (8)  $\beta_i(t - 1) = \sum_{j=1}^{\sigma} p(s^i \xrightarrow{w_{t-1}} s^j) \beta_j(t)$

**D. Training a simple Markov chain**

- (9) A simple Markov chain



$$(10) \quad p_e(s^i \xrightarrow{w^k} s^j) = \frac{|s^i \xrightarrow{w^k} s^j|}{\sum_{t=1, m=1}^{\sigma, \omega} |s^i \xrightarrow{w^m} s^t|}$$

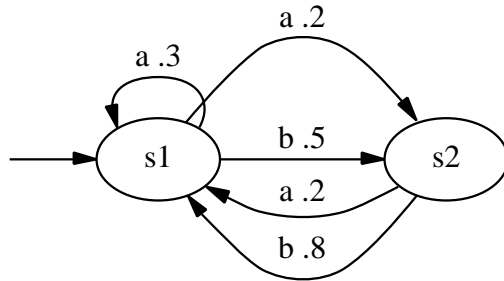
$$(11) \quad s^1 \xrightarrow{a} s^1 \xrightarrow{b} s^2 \xrightarrow{a} s^1 \xrightarrow{b} s^2 \xrightarrow{a} s^1 \xrightarrow{b} s^2 \xrightarrow{b} s^1$$

|      | arc                       | count | probability |
|------|---------------------------|-------|-------------|
| (12) | $s^1 \xrightarrow{a} s^1$ | 1     | .25         |
|      | $s^1 \xrightarrow{b} s^2$ | 3     | .75         |
|      | $s^2 \xrightarrow{a} s^1$ | 2     | .66         |
|      | $s^2 \xrightarrow{b} s^1$ | 1     | .33         |

### E. Baum-Welch

$$(13) \quad |s^i \xrightarrow{w^k} s^j| = \frac{1}{p(w_{1,n})} \sum_{t=1}^n \alpha_i(t) p(s^i \xrightarrow{w^k} s^j) \beta_j(t+1)$$

(14) A simple Hidden Markov Model



(15) We will train the HMM above with *aabb*. We first calculate the values for  $\alpha$  and  $\beta$ .

|            | $\alpha(1)$ | $\alpha(2)$ | $\alpha(3)$ | $\alpha(4)$ | $\alpha(5)$ |
|------------|-------------|-------------|-------------|-------------|-------------|
| (16) $s^1$ | 1           | .3          | .13         | .048        | .052        |
| $s^2$      | 0           | .2          | .06         | .065        | .024        |

|            | $\beta(1)$ | $\beta(2)$ | $\beta(3)$ | $\beta(4)$ | $\beta(5)$ |
|------------|------------|------------|------------|------------|------------|
| (17) $s^1$ | .076       | .2         | .4         | .5         | 1          |
| $s^2$      | 0          | .08        | .4         | .8         | 1          |

$$\begin{aligned}
& |s^i \xrightarrow{w^k} s^j| = \frac{1}{p(w_{1,n})} \sum_{t=1}^n \alpha_i(t) p(s^i \xrightarrow{w^k} s^j) \beta_j(t+1) \\
(18) \quad & |s^1 \xrightarrow{a} s^1| = \frac{1}{.074} [(1 \times .3 \times .2) + (.3 \times .3 \times .4)] = 1.2972 \\
& |s^1 \xrightarrow{a} s^2| = \frac{1}{.074} [(1 \times .2 \times .08) + (.3 \times .2 \times .4)] = .5405 \\
& |s^1 \xrightarrow{b} s^2| = \frac{1}{.074} [(.13 \times .5 \times .8) + (.048 \times .5 \times 1)] = 1.027 \\
& |s^2 \xrightarrow{a} s^1| = \frac{1}{.074} [(0 \times .2 \times .2) + (.2 \times .2 \times .4)] = .2162 \\
& |s^2 \xrightarrow{b} s^1| = \frac{1}{.074} [(.06 \times .8 \times .5) + (.065 \times .8 \times 1)] = 1.027
\end{aligned}$$

$$\begin{aligned}
(19) \quad & \begin{array}{|l|l|} \hline |s^i \xrightarrow{w^k} s^j| & p_e \\ \hline |s^1 \xrightarrow{a} s^1| & .45 \\ |s^1 \xrightarrow{a} s^2| & .19 \\ |s^1 \xrightarrow{b} s^2| & .36 \\ |s^2 \xrightarrow{a} s^1| & .17 \\ |s^2 \xrightarrow{b} s^1| & .83 \\ \hline \end{array}
\end{aligned}$$

## References

- CHARNIAK, EUGENE. 1993. *Statistical Language Learning*. Cambridge: MIT Press.
- HAMMOND, MICHAEL, 2003. *Statistical natural language processing*. U. of Arizona.
- HOPCROFT, J.E., & J.D. ULLMAN. 1979. *Introduction to Automata Theory, Languages, and Computation*. Reading: Addison-Wesley.