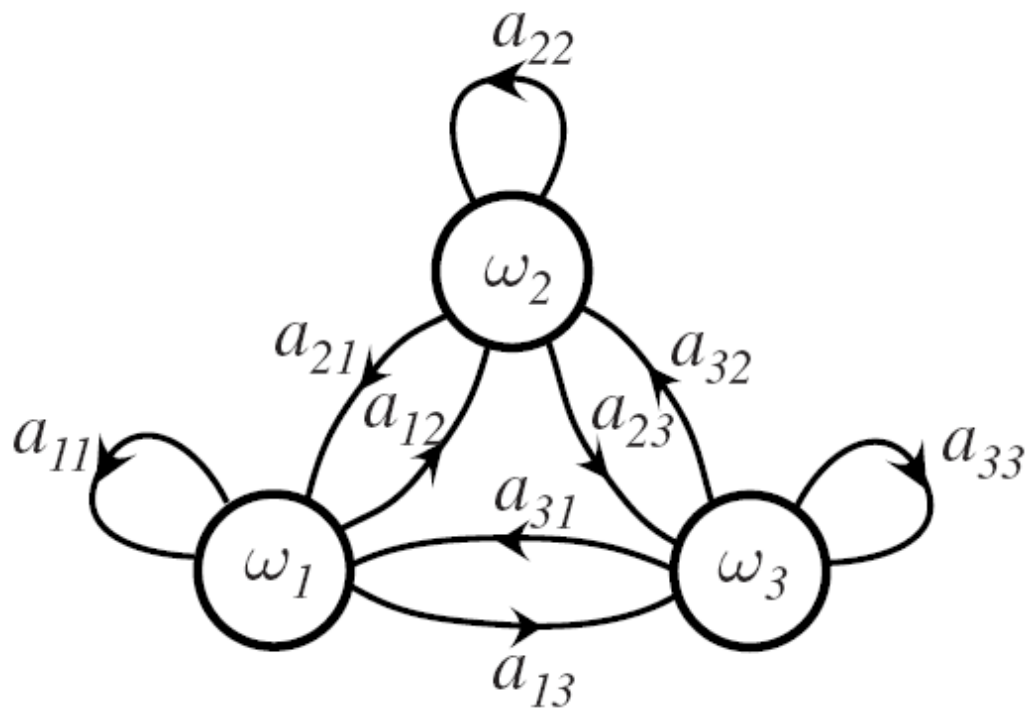# Ch3-p6: Hidden Markov Models

- Markov Chains

  - Goal: make a sequence of decisions

    - Processes that unfold in time, states at time $t$ are influenced by a state at time $t$-1

    - Applications: speech recognition, gesture recognition, parts of speech tagging and DNA sequencing,

    - Any temporal process without memory
      $\omega^T = \{\omega(1), \omega(2), \omega(3), \dots, \omega(T)\}$ sequence of states
      We might have $\omega^6 = \{\omega_1, \omega_4, \omega_2, \omega_2, \omega_1, \omega_4\}$

    - The system can revisit a state at different steps and not every state need to be visited

- First-order Markov models
  - Our productions of any sequence is described by the transition probabilities; the time-independent probability of having state $\omega_j$ at step $t + 1$ given that the state at time $t$ was $\omega_i$

$$P(\omega_j(t + 1) \mid \omega_i(t)) = a_{ij} \quad (a_{ij} \neq a_{ji}, \text{ in general})$$

Figure 3.9: The discrete states, $\omega_i$, in a basic Markov model are represented by nodes, and the transition probabilities, $a_{ij}$, by links. In a first-order discrete time Markov model, at any step $t$ the full system is in a particular state $\omega(t)$. The state at step $t+1$ is a random function that depends solely on the state at step $t$ and the transition probabilities.

We are given a particular model $\theta = (a_{ij}, \omega^T) \rightarrow$
$P(\omega^6 | \theta) = a_{14} \cdot a_{42} \cdot a_{22} \cdot a_{21} \cdot a_{14} \cdot P(\omega(1) = \omega_1)$

$\omega^6 = \{\omega_1, \omega_4, \omega_2, \omega_2, \omega_1, \omega_4\}$

Example: speech recognition

"production of spoken words"
Production of the word: "pattern" represented by phonemes
/p/ /a/ /tt/ /er/ /n/ // ( // = silent state)
Transitions from /p/ to /a/, /a/ to /tt/, /tt/ to er/, /er/ to /n/ and /n/ to a silent state
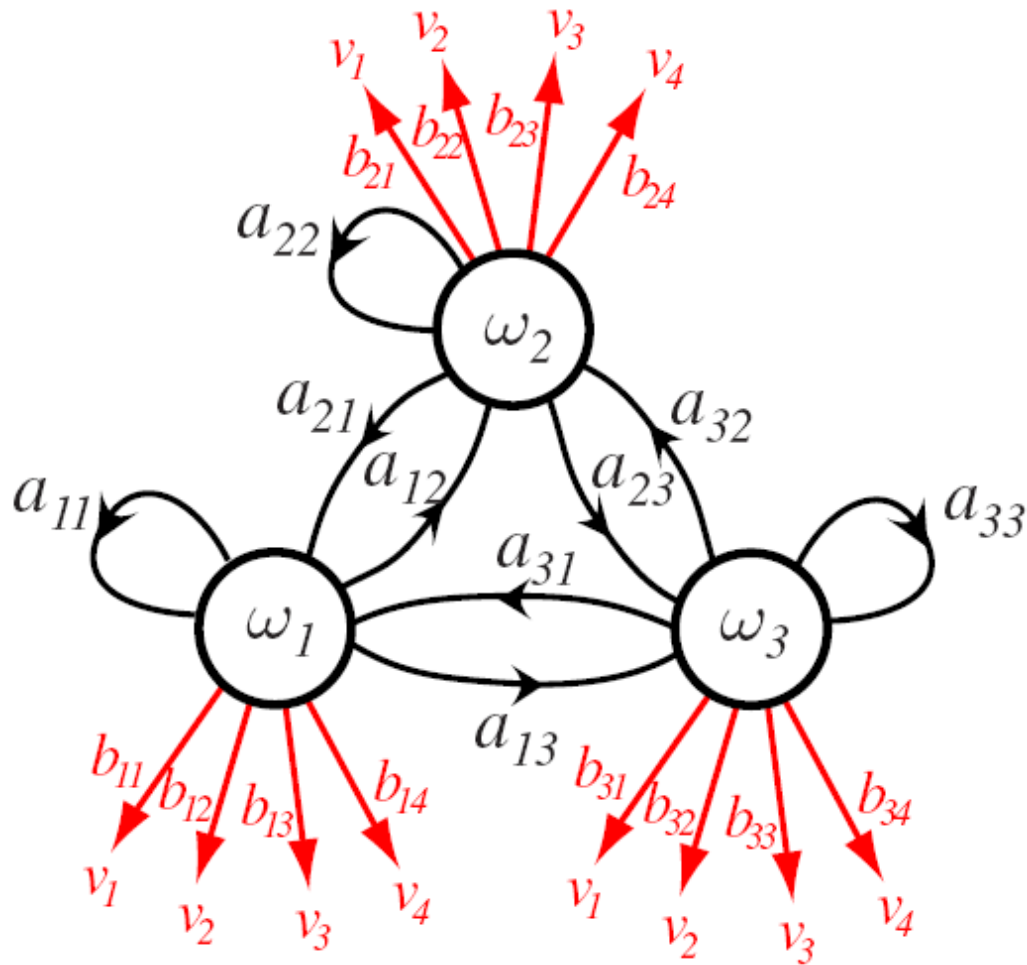
– *visible states: states* which are directly accessible to external measurement

# Hidden Markov Model: Extension of Markov Chains

- We define particular sequence of such visible states as $\mathbf{V}^T = \{v(1), v(2), ..., v(T)\}$ and thus we might have $\mathbf{V}^6 = \{v_5, v_1, v_1, v_5, v_2, v_3\}$.

Because we have access only to the visible states, while the $\omega_i$ are unobservable, such a full model is called a *Hidden Markov Model.*

Thus, a hidden Markov model is a standard Markov process augmented by a set of measurable states and several probabilistic relations between those states and the hidden states.

Figure 3.10: Three hidden units in an HMM and the transitions between them are shown in black while the visible states and the emission probabilities of visible states are shown in red. This model shows all transitions as being possible; in other HMMs, some such candidate transitions are not allowed.

# Hidden Markov Model Computation

- Markov networks are strictly *causal* — the probabilities depend only upon *previous* states. A Markov model is called *ergodic* if every one of the states has a non-zero probability of occurring given some starting state.

- A final or absorbing state $\omega_0$ is one which, if entered, is never left (i.e., $a_{00} = 1$).

$$
\begin{aligned}
a_{ij} &= P(\omega_j(t+1)|\omega_i(t)) \\
b_{jk} &= P(v_k(t)|\omega_j(t)).
\end{aligned}
$$

$$\sum_j a_{ij} = 1 \quad \text{for all } i \text{ and}$$

$$\sum_k b_{jk} = 1 \quad \text{for all } j,$$

where the limits on the summations are over all hidden states and all visible symbols, respectively.

- 3 problems are associated with this model

  - The evaluation problem

  - The decoding problem

  - The learning problem

- **The Evaluation problem.** Suppose we have an HMM, complete with transition probabilities $a_{ij}$ and $b_{jk}$. Determine the probability that a particular sequence of visible states $\mathbf{V}^T$ was generated by that model.

- **The Decoding problem.** Suppose we have an HMM as well as a set of observations $\mathbf{V}^T$. Determine the most likely sequence of *hidden* states $\boldsymbol{\omega}^T$ that led to those observations.

- **The Learning problem.** Suppose we are given the coarse structure of a model (the number of states and the number of visible states) but *not* the probabilities $a_{ij}$ and $b_{jk}$. Given a set of training observations of visible symbols, determine these parameters.

- **The evaluation problem**

  The probability that the model produces a sequence $\mathbf{V}^T$ of visible states is:

  $$P(\mathbf{V}^T \mid \boldsymbol{\Theta}) = \sum_{r=1}^{r_{max}} P(\mathbf{V}^T \mid \boldsymbol{\omega}_r^T) P(\boldsymbol{\omega}_r^T) \; ; \quad r_{max} = c^T$$

  parameters

  $c$ # of hidden states

where each $r$ indexes a particular sequence

$$\boldsymbol{\omega}_r^T = \left\{ \omega(1), \; \omega(2), \; ... \, , \omega(T) \right\} \; \text{of } T \text{ hidden states}$$

(1) $\quad P(\mathbf{V}^T \mid \boldsymbol{\omega}_r^T) = \prod_{t=1}^{t=T} P(v(t) \mid \omega(t))$   Conditional Independence

(2) $\quad P(\boldsymbol{\omega}_r^T) = \prod_{t=1}^{t=T} P(\omega(t) \mid \omega(t-1))$   Makov chain of order 1

Using equations (1) and (2), we can write:

$$P(\mathbf{V}^T \mid \mathbf{\Theta}) = \sum_{r=1}^{r_{max}} \prod_{t=1}^{t=T} P(v(t) \mid \omega(t)) \, P(\omega(t) \mid \omega(t-1))$$
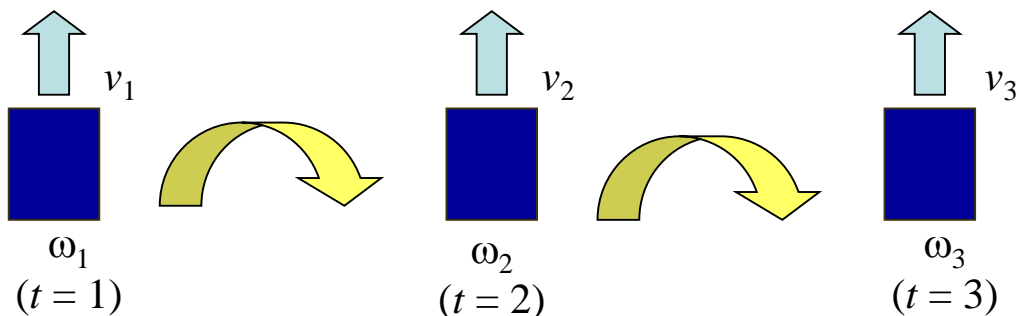
Interpretation: The probability that we observe the particular sequence of $T$ visible states $\mathbf{V}^T$ is equal to the sum over all $r_{max}$ possible sequences of hidden states of the conditional probability that the system has made a particular transition multiplied by the probability that it then emitted the visible symbol in our target sequence.

This is an $O(c^T T)$ calculation, which is quite prohibitive in practice. For instance, if $c = 10$ and $T = 20$, we must perform on the order of $10^{21}$ calculations!!!
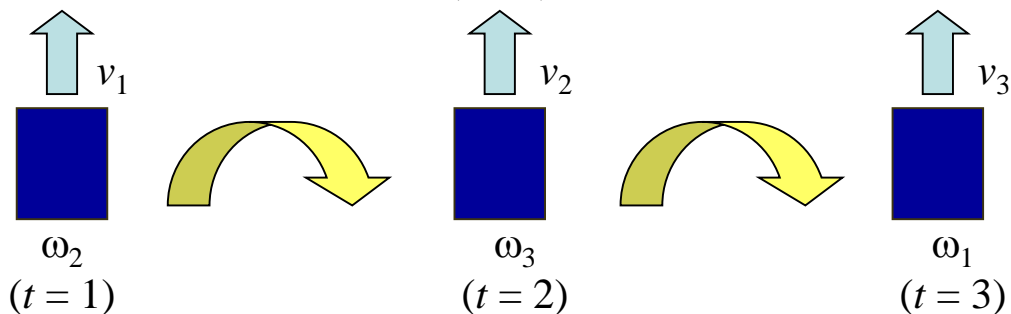
**Example:** Let $\omega_1$, $\omega_2$, $\omega_3$ be the hidden states; $v_1$, $v_2$, $v_3$ be the visible states and $\mathbf{V}^3 = \{v_1, v_2, v_3\}$ is the sequence of visible states

$P(\{v_1, v_2, v_3\}|\Theta) = P(\omega_1).P(v_1 \mid \omega_1).P(\omega_2 \mid \omega_1).P(v_2 \mid \omega_2).P(\omega_3 \mid \omega_2).P(v_3 \mid \omega_3)$
$+P(\omega_2).P(v_1 \mid \omega_2).P(\omega_3 \mid \omega_2).P(v_2 \mid \omega_3).P(\omega_1 \mid \omega_3).P(v_3 \mid \omega_1)+\ldots+$ (possible terms in the sum = all possible ($3^3$= 27) cases !)
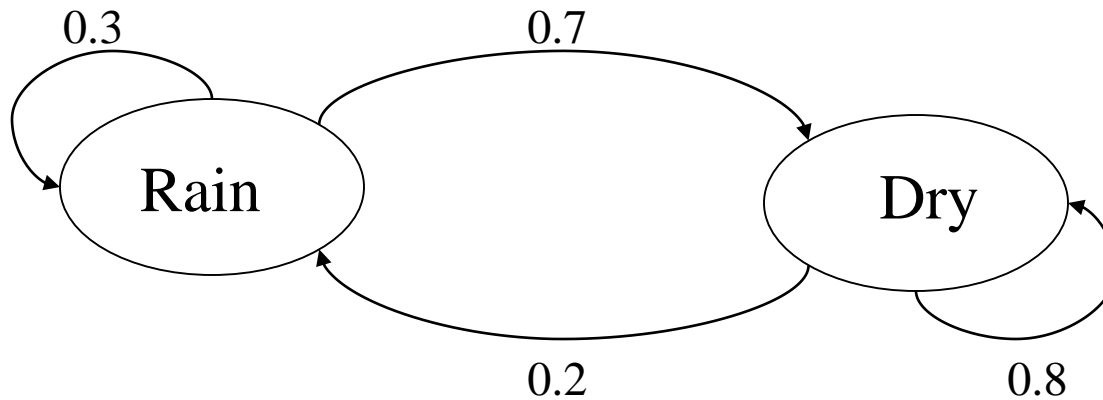
*First possibility:*

$v_1$ $\qquad$ $v_2$ $\qquad$ $v_3$

$\omega_1$ $\qquad$ $\omega_2$ $\qquad$ $\omega_3$
$(t = 1)$ $\qquad$ $(t = 2)$ $\qquad$ $(t = 3)$

*Second Possibility:*

$v_1$ $\qquad$ $v_2$ $\qquad$ $v_3$

$\omega_2$ $\qquad$ $\omega_3$ $\qquad$ $\omega_1$
$(t = 1)$ $\qquad$ $(t = 2)$ $\qquad$ $(t = 3)$

Therefore:

$$P(\{v_1, v_2, v_3\} \mid \Theta) = \sum_{\substack{\text{possible sequence} \\ \text{of hidden states}}} \prod_{t=1}^{t=3} P(v(t) \mid \omega(t)).P(\omega(t) \mid \omega(t-1))$$

12

# Example of Markov Model



- Two states : 'Rain' and 'Dry'.
- Transition probabilities: $P(\text{'Rain'}|\text{'Rain'})=0.3$ , $P(\text{'Dry'}|\text{'Rain'})=0.7$ , $P(\text{'Rain'}|\text{'Dry'})=0.2$, $P(\text{'Dry'}|\text{'Dry'})=0.8$
- Initial probabilities: say $P(\text{'Rain'})=0.4$ , $P(\text{'Dry'})=0.6$ .
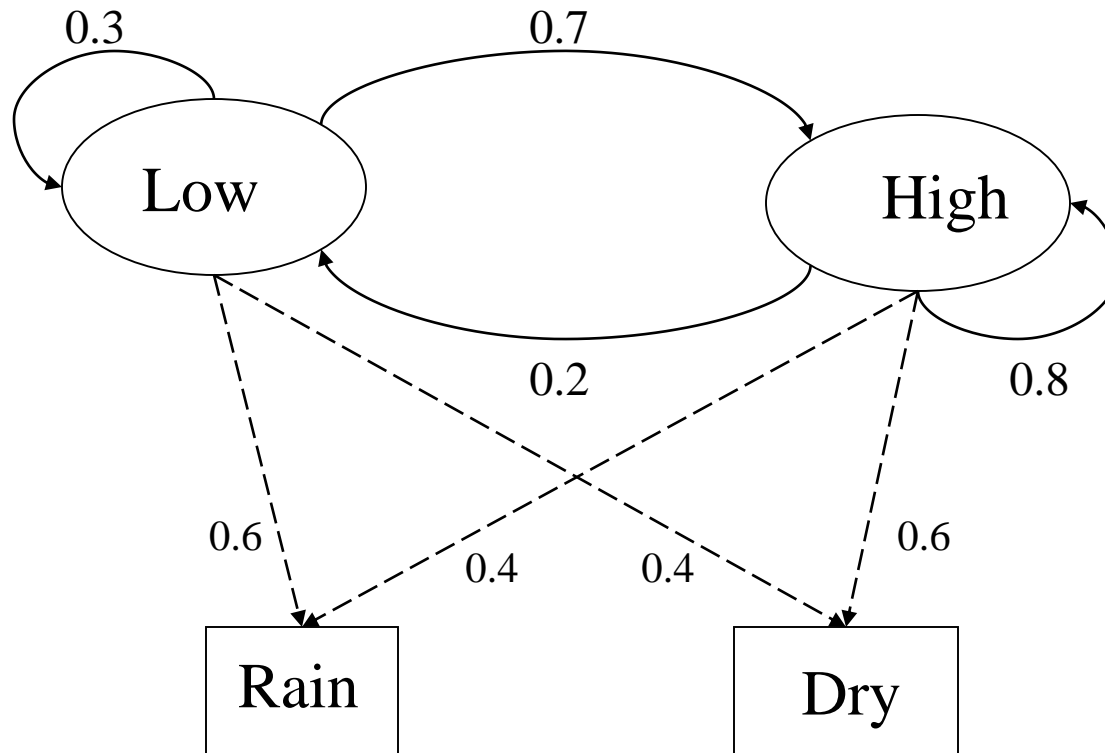
# Calculation of sequence probability

• By Markov chain property, probability of state sequence can be found by the formula:

$$P(s_{i1}, s_{i2}, \ldots, s_{ik}) = P(s_{ik} \mid s_{i1}, s_{i2}, \ldots, s_{ik-1})P(s_{i1}, s_{i2}, \ldots, s_{ik-1})$$
$$= P(s_{ik} \mid s_{ik-1})P(s_{i1}, s_{i2}, \ldots, s_{ik-1}) = \ldots$$
$$= P(s_{ik} \mid s_{ik-1})P(s_{ik-1} \mid s_{ik-2}) \ldots P(s_{i2} \mid s_{i1})P(s_{i1})$$

• Suppose we want to calculate a probability of a sequence of states in our example,  {'Dry', 'Dry', 'Rain', 'Rain'}.

$$\text{P}(\{\text{'Dry', 'Dry', 'Rain', 'Rain'}\}) =$$
$$\text{P}(\text{'Rain'}|\text{'Rain'}) \, \text{P}(\text{'Rain'}|\text{'Dry'}) \, \text{P}(\text{'Dry'}|\text{'Dry'}) \, \text{P}(\text{'Dry'}) =$$
$$= 0.3 \times 0.2 \times 0.8 \times 0.6 = 0.0288$$

# Example of Hidden Markov Model

# Example of Hidden Markov Model

• Two states : 'Low' and 'High' atmospheric pressure.

• Two observations : 'Rain' and 'Dry'.

• Transition probabilities: $P('Low'|'Low')=0.3$ ,

$P('High'|'Low')=0.7$ , $P('Low'|'High')=0.2$,

$P('High'|'High')=0.8$

• Observation probabilities : $P('Rain'|'Low')=0.6$ ,

$P('Dry'|'Low')=0.4$ , $P('Rain'|'High')=0.4$ ,

$P('Dry'|'High')=0.3$ .

• Initial probabilities: say $P('Low')=0.4$ , $P('High')=0.6$ .

# Calculation of observation sequence probability

• Suppose we want to calculate a probability of a sequence of observations in our example, {'Dry','Rain'}.

• Consider all possible hidden state sequences:

$P(\{\text{'Dry','Rain'}\}) = P(\{\text{'Dry','Rain'}\}, \{\text{'Low','Low'}\}) + P(\{\text{'Dry','Rain'}\}, \{\text{'Low','High'}\}) + P(\{\text{'Dry','Rain'}\}, \{\text{'High','Low'}\}) + P(\{\text{'Dry','Rain'}\}, \{\text{'High','High'}\})$

where first term is :

$P(\{\text{'Dry','Rain'}\}, \{\text{'Low','Low'}\}) = P(\{\text{'Dry','Rain'}\} | \{\text{'Low','Low'}\}) \; P(\{\text{'Low','Low'}\}) = P(\text{'Dry'}|\text{'Low'})P(\text{'Rain'}|\text{'Low'}) \; P(\text{'Low'})P(\text{'Low'}|\text{'Low'})$
$= 0.4 \times 0.4 \times 0.6 \times 0.4 \times 0.3 = 0.01152$

# Efficient Calculations: Forward-Backward algorithms

A computationally simpler algorithm for the same goal is as follows.

We can calculate $P(\mathbf{V}^T)$ recursively, since each term $P(v(t)|\omega(t))P(\omega(t)|\omega(t-1))$ involves only $v(t)$, $\omega(t)$ and $\omega(t-1)$. We do this by defining <span style="color:red">forward variable</span> $\alpha_j(t)$

$$\alpha_j(t) = \begin{cases} 0 & t = 0 \text{ and } j \neq \text{initial state} \\ 1 & t = 0 \text{ and } j = \text{initial state} \\ \left[\sum_i \alpha_i(t-1)a_{ij}\right] b_{jk}(v(t)) & \text{otherwise} \end{cases}$$

$b_{jk}(v(t))$ means the transition probability $b_{jk}$ selected by the visible state emitted at time $t$. Thus the only non-zero contribution to the sum is for the index $k$ which matches the visible state $v(t)$.

Thus $\alpha_j(t)$ represents the probability that our HMM is in hidden state $\omega_j$ at step $t$ having generated the first $t$ elements of $\mathbf{V}^T$ i.e. $\alpha_j(t) = P(v(1), v(2), ..., v(t), \omega_j(t) | \Theta)$

**Algorithm 2 (HMM Forward)**

1 <u>**initialize**</u> $\omega(1), t = 0, a_{ij}, b_{jk}$, visible sequence $\mathbf{V}^T, \alpha(0) = 1$
2   <u>**for**</u> $t \leftarrow t + 1$
3       $\alpha_j(t) \leftarrow \sum_{i=1}^{c} \alpha_i(t-1) a_{ij} b_{jk}(v(t))$
4   <u>**until**</u> $t = T$
5 <u>**return**</u> $P(\mathbf{V}^T) \leftarrow \alpha_0(T)$
6 <u>**end**</u>

$\alpha_0$ denotes the probability of the associated sequence ending to the known final state.
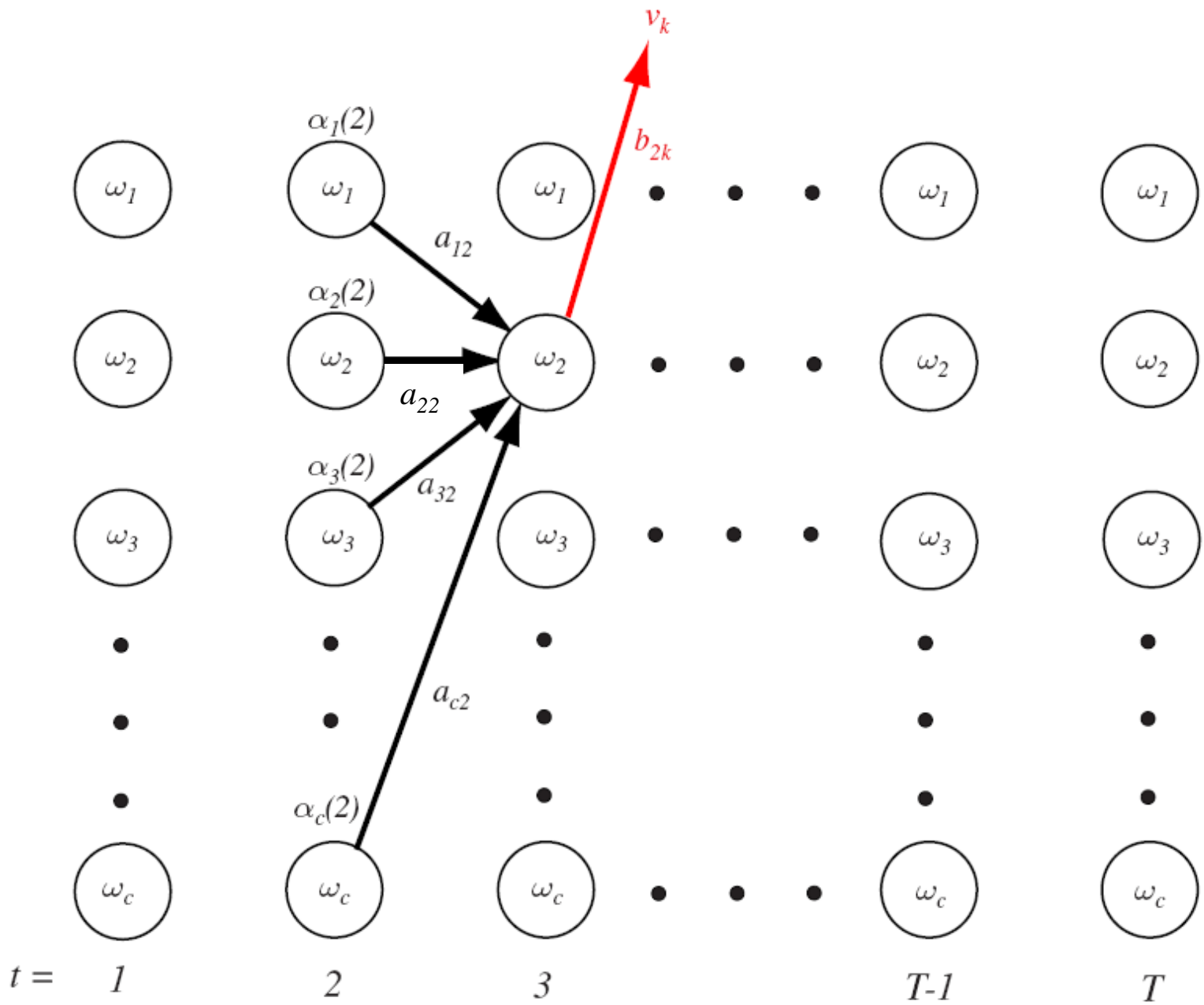*It is $O(c^2 T)$. For $c = 10$, $T = 20 \rightarrow 2000$ calculations*

## Algorithm 3 (HMM Backward)

1 __initialize__ $\omega(T), t = T, a_{ij}, b_{jk}$, visible sequence $V^T$
2    __for__ $t \leftarrow t - 1;$

4       $\beta_i(t) \leftarrow \sum_{j=1}^{C} \beta_j(t+1) a_{ij} b_{jk} \left(v(t+1)\right)$

5    __until__ $t = 1$
7 __return__ $P(V^T) \leftarrow \beta_i(0)$ for the known initial state
8 __end__

where

$$\beta_i(t) = \begin{cases} 0 & \omega_i(t) \neq \text{sequence's final state and } t = T \\ 1 & \omega_i(t) = \text{sequence's final state and } t = T \\ \sum_j a_{ij} b_{jk} v(t+1) \beta_j(t+1) & \text{otherwise,} \end{cases}$$

$\beta_i(t) = P$(model will gen. the seq. from $t+1$ to $T$ given $\omega_i(t)$)
$= P(v(t+1), v(t+2), ..., v(T) \mid \omega_i(t), \Theta)$

**Figure 3.11:** The computation of probabilities by the Forward algorithm can be visualized by means of a trellis — a sort of "unfolding" of the HMM through time. Suppose we seek the probability that the HMM was in state $\omega_2$ at $t = 3$ and generated the observed visible up through that step (including the observed visible symbol $v_k$). The probability the HMM was in state $\omega_j$ ($t = 2$) and generated the observed sequence through $t = 2$ is $\alpha_j(2)$ for $j = 1, 2, ..., c$. To find $\alpha_2(3)$ we must sum these and multiply the probability that state $\omega_2$ emitted the observed symbol $v_k$. Formally, for this particular illustration we have

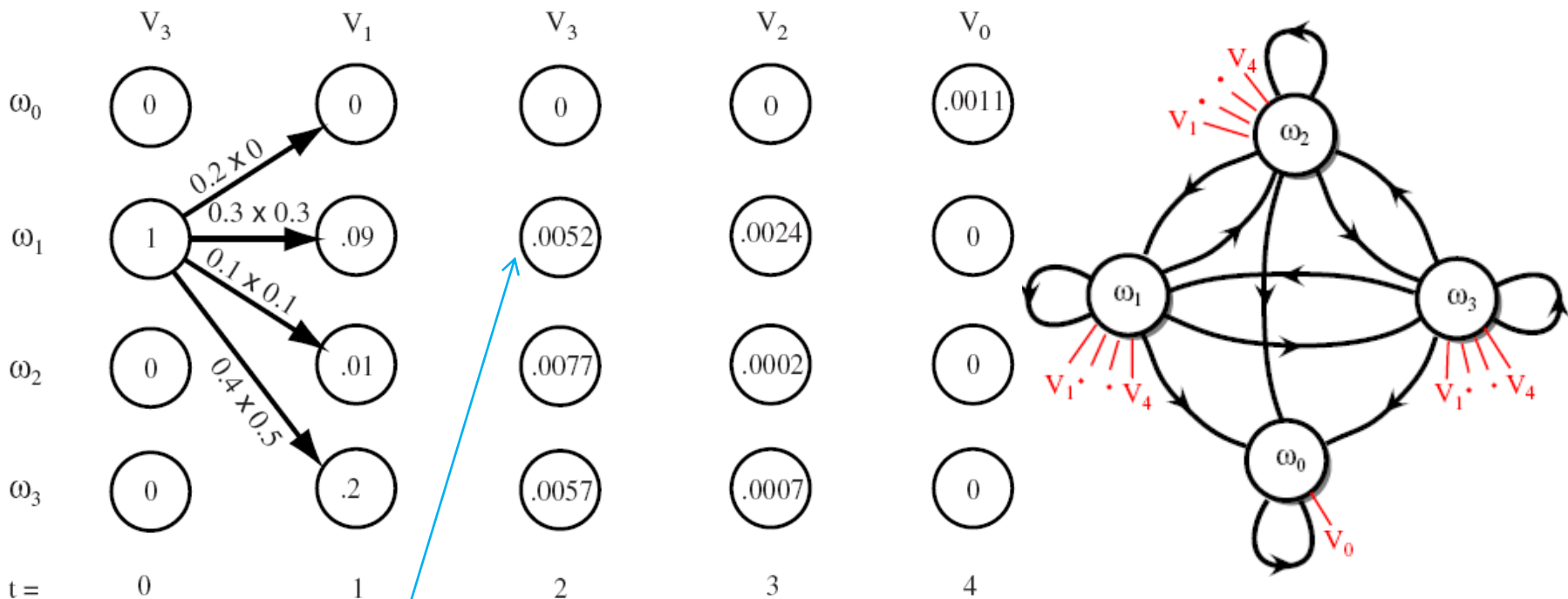$$\alpha_2(3) = b_{2k} \sum_{j=1}^{c} \alpha_j(2) a_{j2}.$$

Consider an HMM such as shown in Fig. 3.10, but with an explicit absorber state and unique null visible symbol $v_0$ with the following transition probabilities (where the matrix indexes begin at 0):

| $a_{ij}$ | $\omega_0$ | $\omega_1$ | $\omega_2$ | $\omega_3$ |
|---|---|---|---|---|
| $\omega_0$ | 1 | 0 | 0 | 0 |
| $\omega_1$ | 0.2 | 0.3 | 0.1 | 0.4 |
| $\omega_2$ | 0.2 | 0.5 | 0.2 | 0.1 |
| $\omega_3$ | 0.8 | 0.1 | 0 | 0.1 |

| $b_{jk}$ | $v_0$ | $v_1$ | $v_2$ | $v_3$ | $v_4$ |
|---|---|---|---|---|---|
| $\omega_0$ | 1 | 0 | 0 | 0 | 0 |
| $\omega_1$ | 0 | 0.3 | 0.4 | 0.1 | 0.2 |
| $\omega_2$ | 0 | 0.1 | 0.1 | 0.7 | 0.1 |
| $\omega_3$ | 0 | 0.5 | 0.2 | 0.1 | 0.2 |

What is the probability it generates the particular sequence $\mathbf{V}^5 = \{v_3, v_1, v_3, v_2, v_0\}$?



$$\alpha_1(2) = b_{13} \sum_{i=0}^{3} \alpha_i(1)a_{i1} = 0.1(0 \times 0 + 0.09 \times 0.3 + 0.01 \times 0.5 + 0.2 \times 0.1) = 0.0052$$

Figure: The HMM (above) consists of four hidden states (one of which is an absorber state, $\omega_0$), each emitting one of five visible states; only the allowable transitions to visible states are shown. The trellis for this HMM is shown below. In each node is $\alpha_i(t)$—the probability the model generated the observed visible sequence up to $t$. For instance, we know that the system was in hidden state $\omega_1$ at $t = 1$, and thus $\alpha_1(0) = 1$ and $\alpha_i(0) = 0$ for $i \neq 1$. The arrows show the calculation of $\alpha_i(1)$. For instance, since visible state $v_1$ was emitted at $t = 1$, we have $\alpha_0(1) = \alpha_1(0)a_{10}b_{10} = 1[0.2 \times 0] = 0$ as shown by the top arrow. Likewise the next highest arrow corresponds to the calculation $\square_1(1) = \alpha_1(0)a_{11}b_{11} = 1[0.3 \times 0.3] = 0.09$. In this example, the calculation of $\alpha_i(1)$ is particularly simple, since only transitions from the known initial hidden state need be considered; all other transitions have zero contribution to $\alpha_i(1)$. For subsequent times, however, the calculation requires a sum over all hidden states at the previous time, as given by line 3 in the Forward algorithm. The probability shown in the final (absorbing) state gives the probability of the full sequence observed, $P(\mathbf{V}^T | \boldsymbol{\theta}) = 0.0011$.
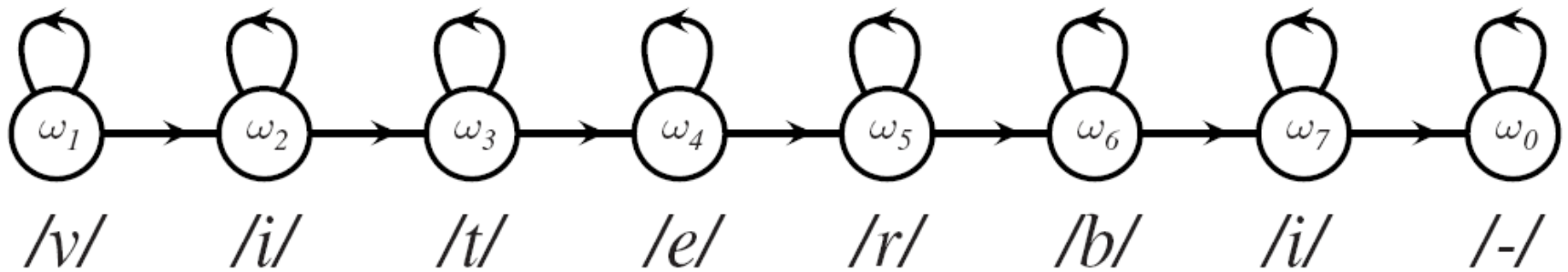
If we denote our model - the $a$'s and $b$'s - by $\boldsymbol{\theta}$, we have by Bayes' formula

$$P(\boldsymbol{\theta}|\mathbf{V}^T) = \frac{P(\mathbf{V}^T|\boldsymbol{\theta})P(\boldsymbol{\theta})}{P(\mathbf{V}^T)}$$

In HMM pattern recognition we would have a number of HMMs, one for each category and classify a test sequence according to the model with the highest probability.

The Forward algorithm gives us $P(\mathbf{V}^T/\boldsymbol{\theta})$.

The prior probability of the model, $P(\boldsymbol{\theta})$, is given by some external source, such as a *language model* in the case of speech.

Figure 3.12: A left-to-right HMM commonly used in speech recognition. For instance, such a model could describe the utterance "viterbi," where $\omega_1$ represents the phoneme /v/, $\omega_2$ represents /i/, ..., and $\omega_0$ a final silent state. Such a left-to-right model is more restrictive than the general HMM in Fig. 3.11, and preclude transitions "back" in time.

# Normalization is Important

- **Normalization is required to avoid such recursive algorithms from accumulating large amounts of computational noise.**

- **We can apply a normalization factor at each step of the calculation:**

$$\alpha'_j(t) = \frac{\alpha_j(t)}{\prod\limits_{i=0}^{t} Q_i}$$

  **where the scale factor, Q, is given by:**

$$Q_i = \sum_{i=1}^{c} \alpha'_i(t) \quad Q_0 = 1$$

- **This is applied once per state per unit time, and simply involves scaling the current $\alpha$'s by their sum at each epoch (e.g., a frame).**

- **Also, likelihoods tend to zero as time increases and can cause underflow. Therefore, it is more common to operate on log probabilities to maintain numerical precision . This converts products to sums but still involves essentially the same algorithm (though an approximation for the log of a sum is used to compute probabilities involving the summations).**

# Classification Using HMMs

If we concatenate our HMM parameters into a single vector, $\mathbf{\theta}$, we can write Bayes formula as:

$$P\left(\mathbf{\theta}\middle|\mathbf{v}^T\right) = \frac{P\left(\mathbf{v}^T\middle|\mathbf{\theta}\right)P(\mathbf{\theta})}{P(\mathbf{v}^T)}$$

The forward algorithm gives us $P\left(\mathbf{v}^T\middle|\mathbf{\theta}\right)$

• We ignore the denominator term (evid.) during the maximization.

• In some applications, we use domain knowledge to compute $P(\mathbf{\theta})$. For example, in speech recognition, this most often represents the probability of a word or sound, which comes from a "language model." It is also possible to use HMMs to model $P(\mathbf{\theta})$ (e.g., statistical language modeling in speech recognition).

• In a typical classification application, there are a set of HMMs, one for each category, and the above calculation is performed for each model ($\theta_i$).

- **The decoding problem (optimal state sequence)**

  Given a sequence of visible states $V^T$, the decoding problem is to find the most probable sequence of hidden states, which is called maximum a posteriori (MAP) estimate of the state sequence.

  This problem can be expressed mathematically as:

  *find the single "best" state sequence (hidden states)*

  $\hat{\omega}(1), \hat{\omega}(2),..., \hat{\omega}(T)$ *such that* :

  $\hat{\omega}(1), \hat{\omega}(2),..., \hat{\omega}(T) =$

  $\underset{\omega(1),\omega(2),...,\omega(T)}{\arg \max} \quad P\big[\omega(1), \omega(2),..., \omega(T), v(1), v(2),..., V(T) \,|\, \Theta\big]$

  $$\mathbf{v}^{T^*} = \arg \underset{r}{ma\, x} P\big(\mathbf{v}^T \big| \boldsymbol{\omega}_r^T\big) P\big(\boldsymbol{\omega}_r^T\big)$$

Note that the summation disappeared, since we want to find only one unique best case !

30

Where: $\lambda = \Theta = [\pi, \mathbf{A}, \mathbf{B}]$ is an HMM

$\pi = P(\omega(1))$ (initial state probability)

$\mathbf{A} = a_{ij} = P(\omega(t+1) = j \mid \omega(t) = i)$

$\mathbf{B} = b_{jk} = P(v(t) = k \mid \omega(t) = j)$

$$\mathbf{A} = \begin{bmatrix} a_{11} & \dots & a_{1c} \\ \vdots & \ddots & \vdots \\ a_{c1} & \dots & a_{cc} \end{bmatrix} \qquad a_{ij} = P\big(\omega_j(t+1) \big| \omega_i(t)\big)$$

$$\mathbf{B} = \begin{bmatrix} b_{11} & \dots & b_{1M} \\ \vdots & \ddots & \vdots \\ \vdots & \vdots & \vdots \\ b_{T1} & & a_{TM} \end{bmatrix} \qquad b_{ij} = P\big(v_j(t) \big| \omega_i(t)\big)$$

Initial state distribution: $\pi(0) = \{\pi_1, \pi_2, \dots, \pi_c\}$ $\quad \pi_i = P(\omega_i \mid t = 0)$

$\pi_i$ is the probability that the model is in state $\omega_i$ at the time $t=0$.
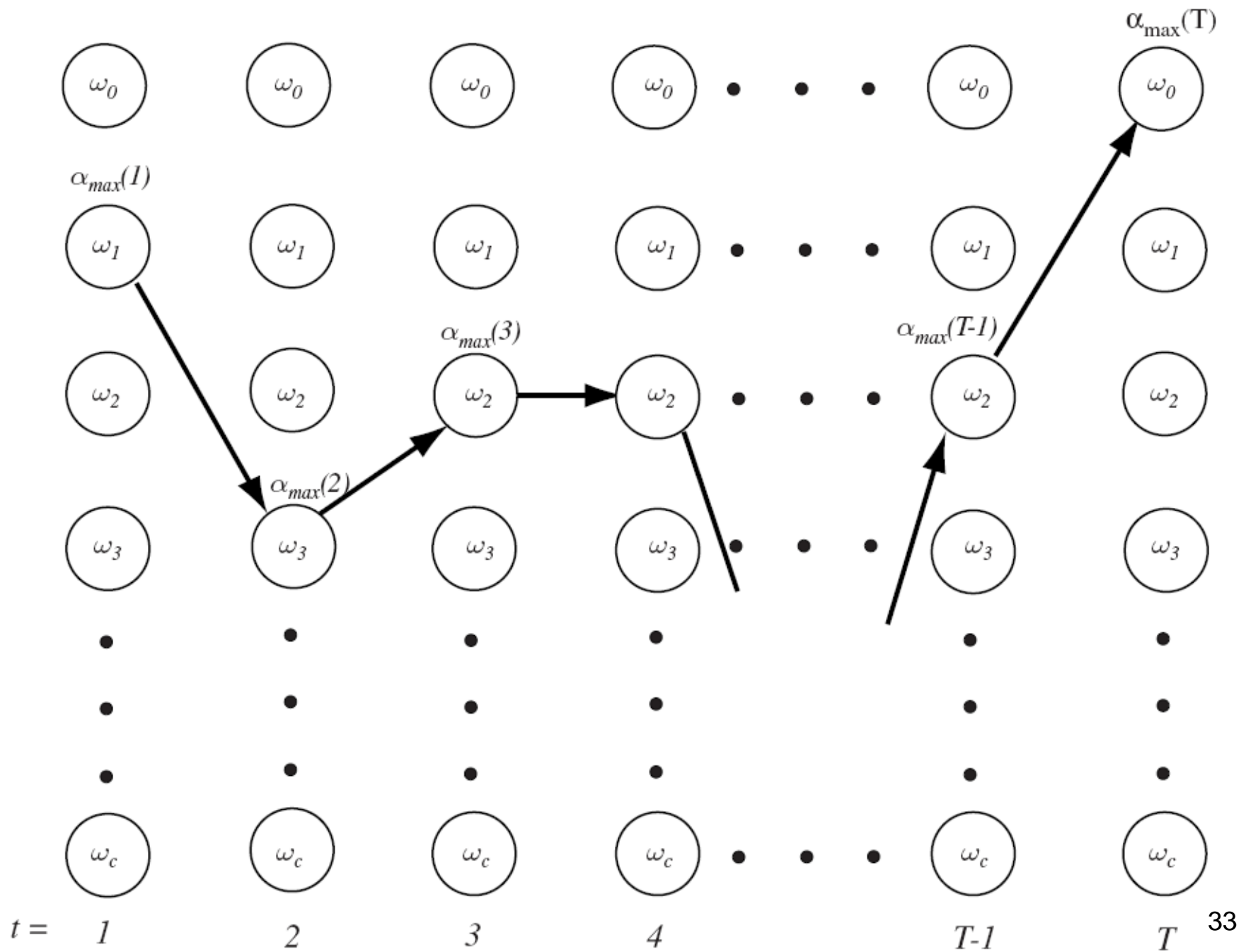
# Algorithm 4: (HMM Decoding)

1 **begin** **initialize** Path $\leftarrow \{\}, t \leftarrow 0$

2        **for** $t \leftarrow t + 1$            Viterbi Algorithm

3            $j \leftarrow j + 1$

4            **for** $j \leftarrow j + 1$

5              $\alpha_j(t) \leftarrow b_{jk}v(t) \sum_{i=1}^{c} \alpha_i(t-1)a_{ij}$

6            **until** $j = c$

7            $j' \leftarrow \arg\max_{j} \alpha_j(t)$

8            $Append$ $\omega_{j'}$ to Path

9        **until** $t = T$

10     **return** Path

11 **end**

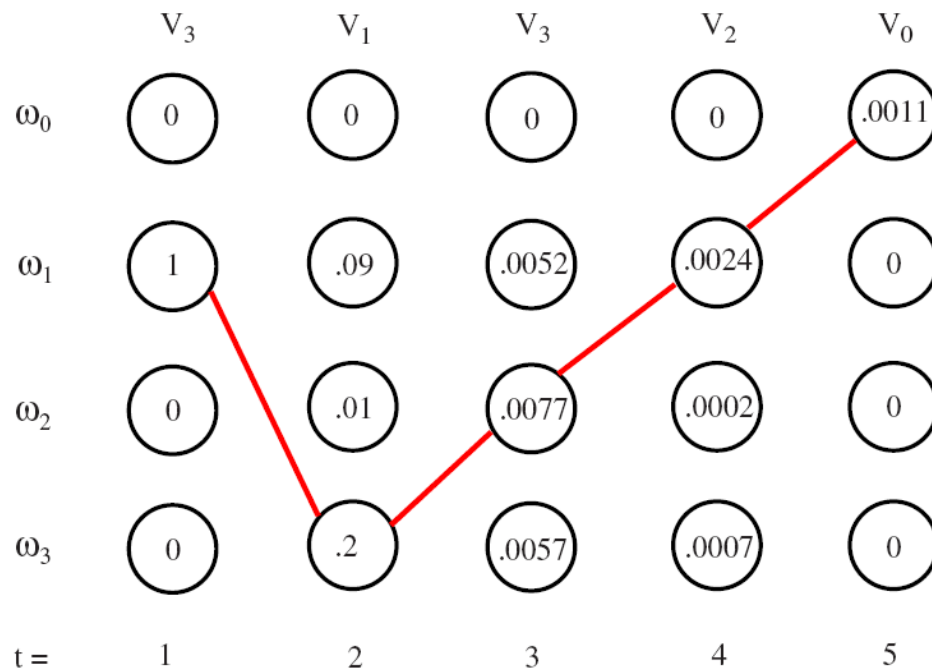$O(c^T T)$ calculation $\rightarrow O(c^2 T)$

**FIGURE 3.13.** Viterbi Decoding Trellis: The decoding algorithm finds at each time step $t$ the state that has the highest probability of having come from the previous step and generated the observed visible state $v_k$. The full path is the sequence of such states. Because this is a local optimization (dependent only upon the single previous time step, not the full sequence), the algorithm does not guarantee that the path is indeed allowable. For instance, it might be possible that the maximum at $t = 5$ is $\omega_1$ and at $t=6$ is $\omega_2$, and thus these would appear in the path. This can even occur if $a_{12} = P(\omega_2 (t + 1)|\omega_1 (t)) = 0$, precluding that transition.

We find the path for the data of Example 4 for the sequence $\{\omega_1, \omega_3, \omega_2, \omega_1, \omega_0\}$. Note especially that the transition from $\omega_3$ to $\omega_2$ is not allowed according to the transition probabilities $a_{ij}$ given in Example 4. The path *locally* optimizes the probability through the trellis.

|  | $V_3$ | $V_1$ | $V_3$ | $V_2$ | $V_0$ |
|---|---|---|---|---|---|
| $\omega_0$ | 0 | 0 | 0 | 0 | .0011 |
| $\omega_1$ | 1 | .09 | .0052 | .0024 | 0 |
| $\omega_2$ | 0 | .01 | .0077 | .0002 | 0 |
| $\omega_3$ | 0 | .2 | .0057 | .0007 | 0 |
| $t =$ | 1 | 2 | 3 | 4 | 5 |

35

- **The learning problem (parameter estimation)**

This third problem consists of determining a method to adjust the model parameters $\Theta = [\pi, \mathbf{A}, \mathbf{B}]$ to satisfy a certain optimization criterion. We need to find the best model

$$\hat{\Theta} = [\hat{\pi}, \hat{\mathbf{A}}, \hat{\mathbf{B}}]$$

Such that to maximize the probability of the observation sequence:

$$\underset{\Theta}{Max}\, P(\mathbf{V}^T \mid \Theta)$$

We use an iterative procedure such as Baum-Welch (Forward-Backward) or Gradient to find this local optimum

- $\alpha_i(t)$ = P(model generates visible sequence up to step $t$ given hidden state $\omega_i(t)$)
- $\beta_i(t)$ = P(model will generate the sequence from $t+1$ to $T$ given $\omega_i(t)$)

$$\beta_i(t) = \begin{cases} 0 & \omega_i(t) \neq \text{sequence's final state and } t = T \\ 1 & \omega_i(t) = \text{sequence's final state and } t = T \\ \sum_j a_{ij} b_{jk} v(t+1) \beta_j(t+1) & \text{otherwise,} \end{cases}$$

By the definition of $\beta_i(T)$ in the above Eq., this will be either 0 (if $\omega_i(T)$ is not the final hidden state) or 1 (if it is). Thus it is clear that $\beta_i(T-1) = \sum_j a_{ij} b_{jk} v(T) \beta_j(T)$.

But the $\alpha_i(t)$ and $\beta_i(t)$ we determined are merely estimates of their true values, since we don't know the actual value of the transition probabilities $a_{ij}$ and $b_{jk}$ in the above Eq.

We now define a posteriori probability $\gamma_i(t) = p\left(\omega_i(t) \mid \mathbf{V}, \boldsymbol{\theta}\right)$

which is the probability of being in state $\omega_i$ at time $t$ for the state sequence $\mathbf{V}$. Note that:

$$\gamma_i(t) = p\left(\omega_i(t) \mid \mathbf{V}, \boldsymbol{\theta}\right) = \frac{p\left(\mathbf{V}, \omega_i(t) \mid \boldsymbol{\theta}\right)}{p\left(\mathbf{V} \mid \boldsymbol{\theta}\right)}$$

$$= \frac{p\left(\mathbf{V}, \omega_i(t) \mid \boldsymbol{\theta}\right)}{\sum_{j=1}^{c} p\left(\mathbf{V}, \omega_j(t) \mid \boldsymbol{\theta}\right)}$$

Also note that because of Markovian conditional independence:

$$\alpha_i(t)\beta_i(t) = p\left(v_1, v_2, \ldots, v_t, \omega_i(t) \mid \boldsymbol{\theta}\right) \times$$

$$p\left(v_{t+1}, v_{t+2}, \ldots, v_T \mid \omega_i(t), \boldsymbol{\theta}\right) = p\left(\mathbf{V}, \omega_i(t) \mid \boldsymbol{\theta}\right)$$

so we can define $\gamma_i(t)$ in term of $\alpha_i(t)$ and $\beta_i(t)$ as:

$$\gamma_i(t) = \frac{\alpha_i(t)\beta_i(t)}{\sum_{j=1}^{c} \alpha_j(t)\beta_j(t)}$$

We also define $\xi_{ij}(t) = p\left(\omega_i(t), \omega_j(t+1) \mid \mathbf{V}, \boldsymbol{\theta}\right)$

which is the probability of being in state $\omega_i$ at time $t$ and being in state $\omega_j$ at time $t+1$ given the model generated the entire training sequence $\mathbf{V}^T$ by any path.

This can also be expanded as:

$$\xi_{ij}(t) = \frac{p\left(\omega_i(t), \omega_j(t+1), \mathbf{V} \mid \boldsymbol{\theta}\right)}{p\left(\mathbf{V} \mid \boldsymbol{\theta}\right)}$$

$$= \frac{\alpha_i(t) a_{ij} b_{jk}(v_k(t+1)) \beta_j(t+1)}{\displaystyle\sum_{j=1}^{c}\sum_{i=1}^{c} \alpha_i(t) a_{ij} b_{jk}(v_k(t+1)) \beta_j(t+1)}$$

or as:

$$\xi_{ij}(t) = \frac{p\left(\omega_i(t) \mid \mathbf{V}\right) p\left(v_{t+1}, v_{t+2}, \ldots, v_T, \omega_j(t+1) \mid \omega_i(t), \boldsymbol{\theta}\right)}{p\left(v_{t+1}, v_{t+2}, \ldots, v_T \mid \omega_i(t), \boldsymbol{\theta}\right)}$$

$$= \frac{\gamma_i(t) a_{ij} b_{jk}(v_k(t+1)) \beta_j(t+1)}{\beta_i(t)} \qquad \gamma_i(t) = \sum_{j=1}^{c} \xi_{ij}(t)$$

If we sum these quantities across time, we can get some useful values i.e., the expression
$$\sum_{t=1}^{T} \gamma_i(t)$$
is the expected number of times in state $\omega_i$ and therefore is the expected number of transitions away from state $\omega_i$. Similarly,
$$\sum_{t=1}^{T-1} \xi_{ij}(t)$$
is the expected number of transitions from state $\omega_i$ to state $\omega_j$ for **V**. These follow from the fact that

$$\sum_{t=1}^{T} \gamma_i(t) = \sum_{t} E\left[I_t(i)\right] = E\left[\sum_{t} I_t(i)\right]$$

$$\sum_{t=1}^{T} \xi_{ij}(t) = \sum_{t} E\left[I_t(i,j)\right] = E\left[\sum_{t} I_t(i,j)\right]$$

Where $I_t(i)$ is an indicator random variable that is 1 when we are in state $\omega_i$ at time $t$ and $I_t(i,j)$ is a random variable that is 1 when we move from state $\omega_i$ to state $\omega_j$ after time $t$.

Thus $\hat{a}_{ij}$ (the estimate of the probability of a transition from $\omega_i(t-1)$ to $\omega_j(t)$) can be found by taking the ratio between the expected number of transitions from $\omega_i$ to $\omega_j$ and the total expected number of any transitions from $\omega_i$.

$$\hat{a}_{ij} = \frac{\sum_{t=1}^{T-1} \xi_{ij}(t)}{\sum_{t=1}^{T-1} \gamma_i(t)}$$

We can obtain an improved estimate $\hat{b}_{ij}$ by calculating the ratio between the frequency that any particular symbol $v_k$ is emitted and that for any symbol.

$$\hat{b}_{ik} = \frac{\text{expected number of times in state } \omega_i \text{ and observe symbol } v_k}{\text{expected number of times in state } \omega_i}$$

$$\hat{b}_{ik} = \frac{\sum\limits_{\substack{t=1 \\ v(t)=v_k}}^{T} \gamma_i(t)}{\sum\limits_{t=1}^{T} \gamma_i(t)}$$

Then, we start with rough or arbitrary estimates of $a_{ij}$ and $b_{ik}$, *calculate* improved estimates by above Eqs., and repeat until some convergence criterion is met.

- **Parameter Updates:**

Baum-Welch Forward-Backward Algorithm

- Parameters Learning Algorithm

Begin initialize

$a_{ij}$, $b_{ik}$, training sequence $V^T$, convergence criterion $\theta$, $z=0$

Do $z \leftarrow z+1$

   compute $\hat{a}(z)$ from $a(z-1)$ and $b(z-1)$

   compute $\hat{b}(z)$ from $a(z-1)$ and $b(z-1)$

   $a_{ij}(z) \leftarrow \hat{a}_{ij}(z)$

   $b_{ik}(z) \leftarrow \hat{b}_{ik}(z)$

Until $max_{i,j,k}$ $\{a_{ij}(z)-a_{ij}(z-1), b_{jk}(z)-b_{ik}(z-1)\} < \theta$

Return $a_{ij} \leftarrow a_{ij}(z); b_{ik} \leftarrow b_{ik}(z)$

End

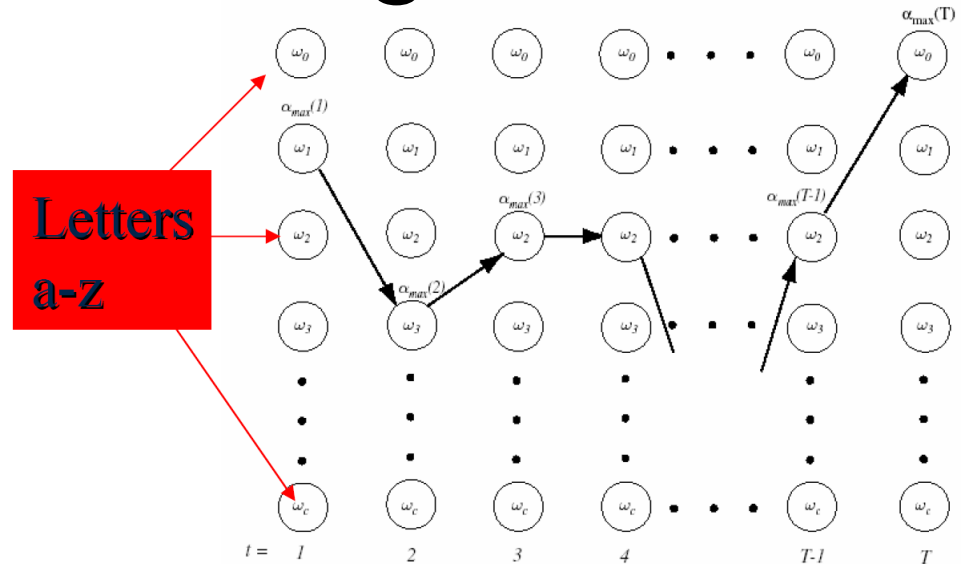$$\hat{a}_{ij} = \frac{\sum_{t=1}^{T-1} \xi_{ij}(t)}{\sum_{t=1}^{T-1} \gamma_i(t)}$$

$$\hat{b}_{ik} = \frac{\sum_{\substack{t=1 \\ v(t)=v_k}}^{T} \gamma_i(t)}{\sum_{t=1}^{T} \gamma_i(t)}$$

# Convergence

- Requires several presentations of each training sequence (fewer than 5 common in speech)

- Another stopping criterion:

    - Overall probability that learning model could have generated the training data

- Seeding the parameters with good initial guesses is very important.

# HMM Word Recognition

- Two approaches:
  - HMM can model all possible words
  - Each state corresponds to each letter of alphabet
  - Letter transition probabilities are calculated for each pair of letters
  - Letter confusion probabilities are symbol probabilities
  - Decoding problem gives most likely word

- Separate HMMs are used to model each word
  - Evaluation problem gives probability of observation which is used as a class-conditional probability



- Each word, e.g., cat, dog, etc, has an associated HMM
- For a test utterance determine which model has highest probability
- HMMs for speech are left-to-right models
- HMM produces a class conditional class-probability

46