

# Hidden Markov Models

Slides by Margrit Betke, Yiwen Gu

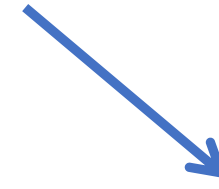
Reading: Rabiner, Proceedings of the IEEE, 77(2), 1989, up to page 275

# Agenda (10/29, 10/31)

- From State Machine To Markov Model
- Working with Hidden Markov Models

Finite State Machine (FSM)

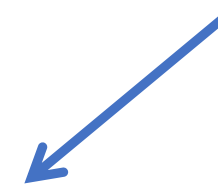
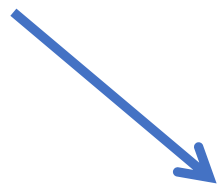
= Automata



Markov Network

= FSM with Transition Probabilities

Finite State Machine with  
Deterministic Outputs

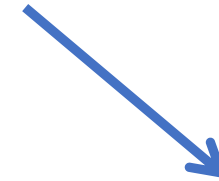


Hidden Markov Model

= Markov Network with Output Probabilities

Finite State Machine (FSM)

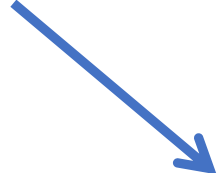
= Automata



Markov Network

= FSM with Transition Probabilities

Finite State Machine with  
Deterministic Outputs

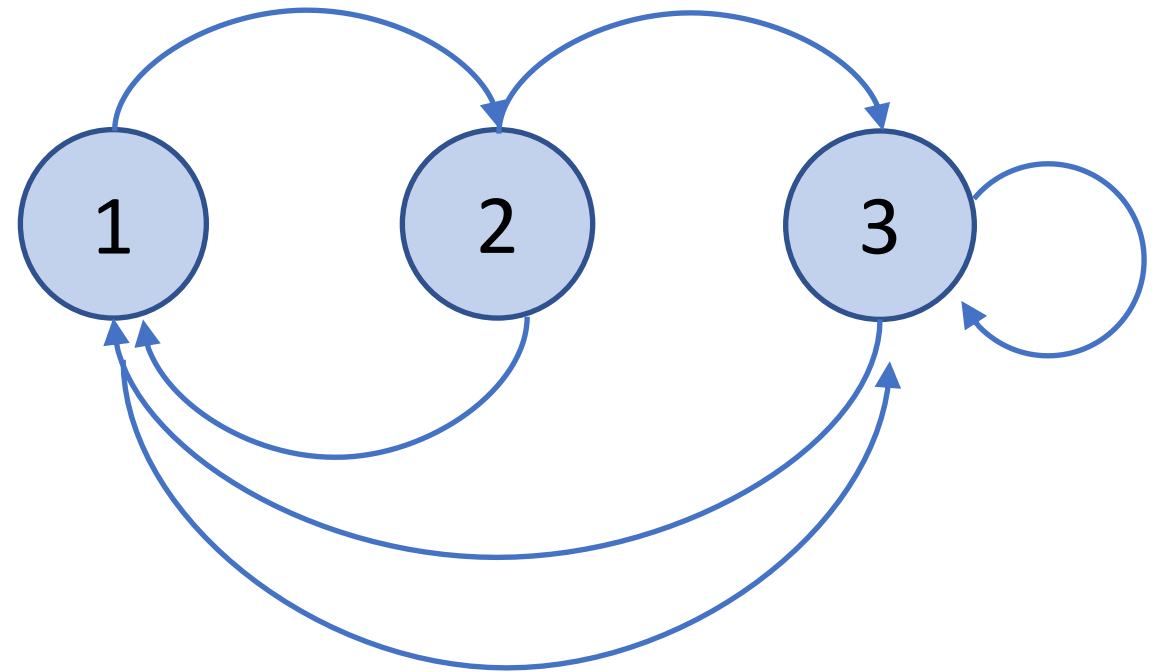


Hidden Markov Model

= Markov Network with Output Probabilities

# Finite State Machine (FSM) = Automata

- A state machine is a machine's AI logic in graph form.
- Key Concepts:
  - States/nodes/vertices
  - Transitions/edges/arcs

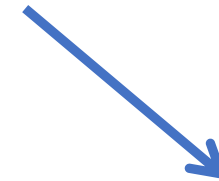


# States =  $N = 3$ ,

Set of states  $Q = \{q_1, \dots, q_N\}$

Finite State Machine (FSM)

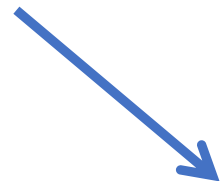
= Automata



Markov Network

= FSM with Transition Probabilities

Finite State Machine with  
Deterministic Outputs

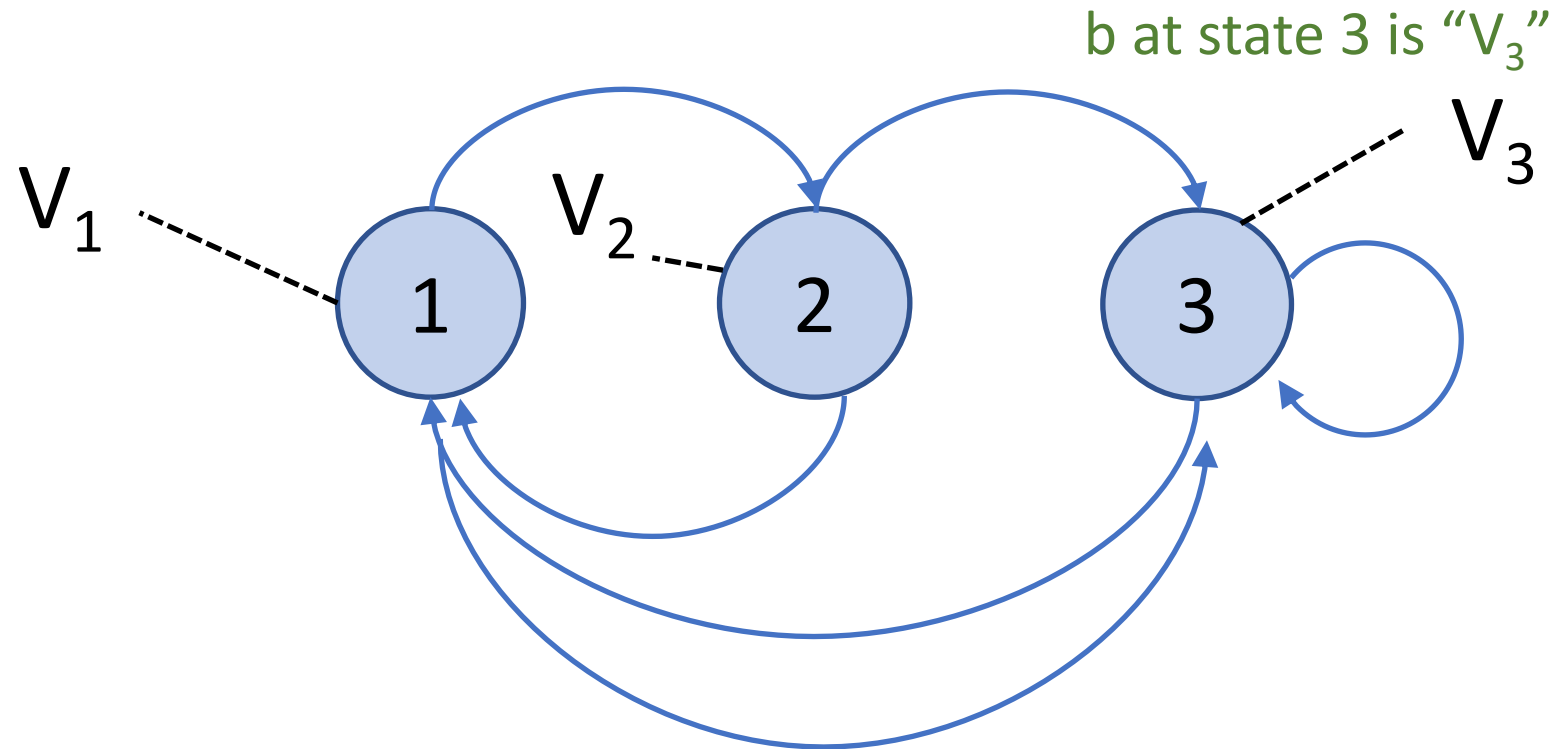


Hidden Markov Model

= Markov Network with Output Probabilities

# FSM with Deterministic Outputs (b's)

- Each state outputs a **symbol** deterministically  
i.e., for every state and input combination, there is a specific, predefined output).
- The symbols in an FSM form a **vocabulary**.
- The set of all possible vocabularies (sequences of the symbols) is called **language**.
- Generator and/or Recognizer



# States =  $N = 3$ , Set of states  $Q = \{q_1, \dots, q_N\}$

# Outputs =  $M = 3$ ,  $\{V_1, V_2, V_3\}$

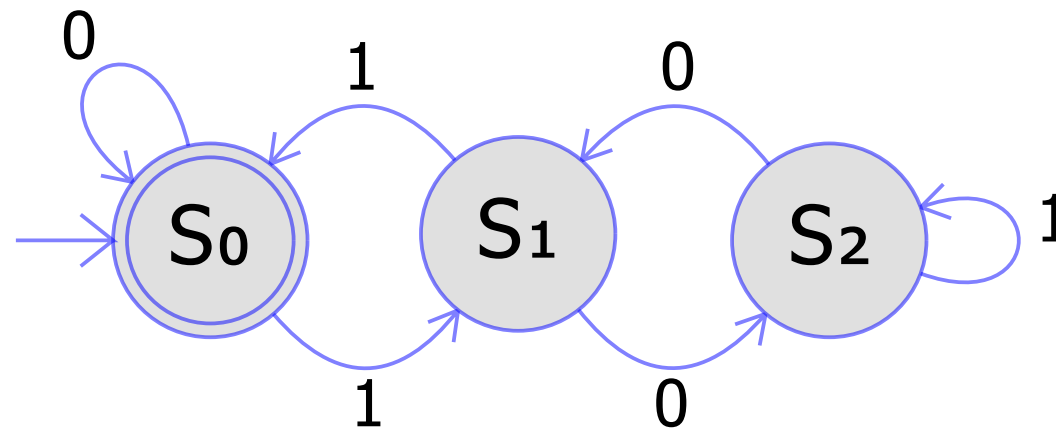
# Example: Deterministic FSM as a Recognizer

## Formal Definition:

A *Deterministic Finite Automaton (Acceptor)* (DFA) is a tuple  $M = \langle Q, \Sigma, \delta, q_0, F \rangle$

- $Q$  - a finite set of **states**
- $\Sigma$  - a finite set of input symbol **alphabet**
- $\delta$  - a transition **function**
- $q_0 \in Q$  - the **start state**
- $F \subseteq Q$  - the **final (or "accepting") states**

Q: What is this doing?

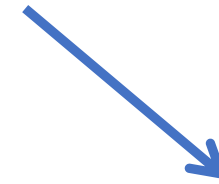


Source : Wiki



Finite State Machine (FSM)

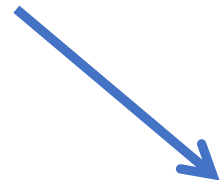
= Automata



Markov Network

= FSM with Transition Probabilities

Finite State Machine with  
Deterministic Outputs



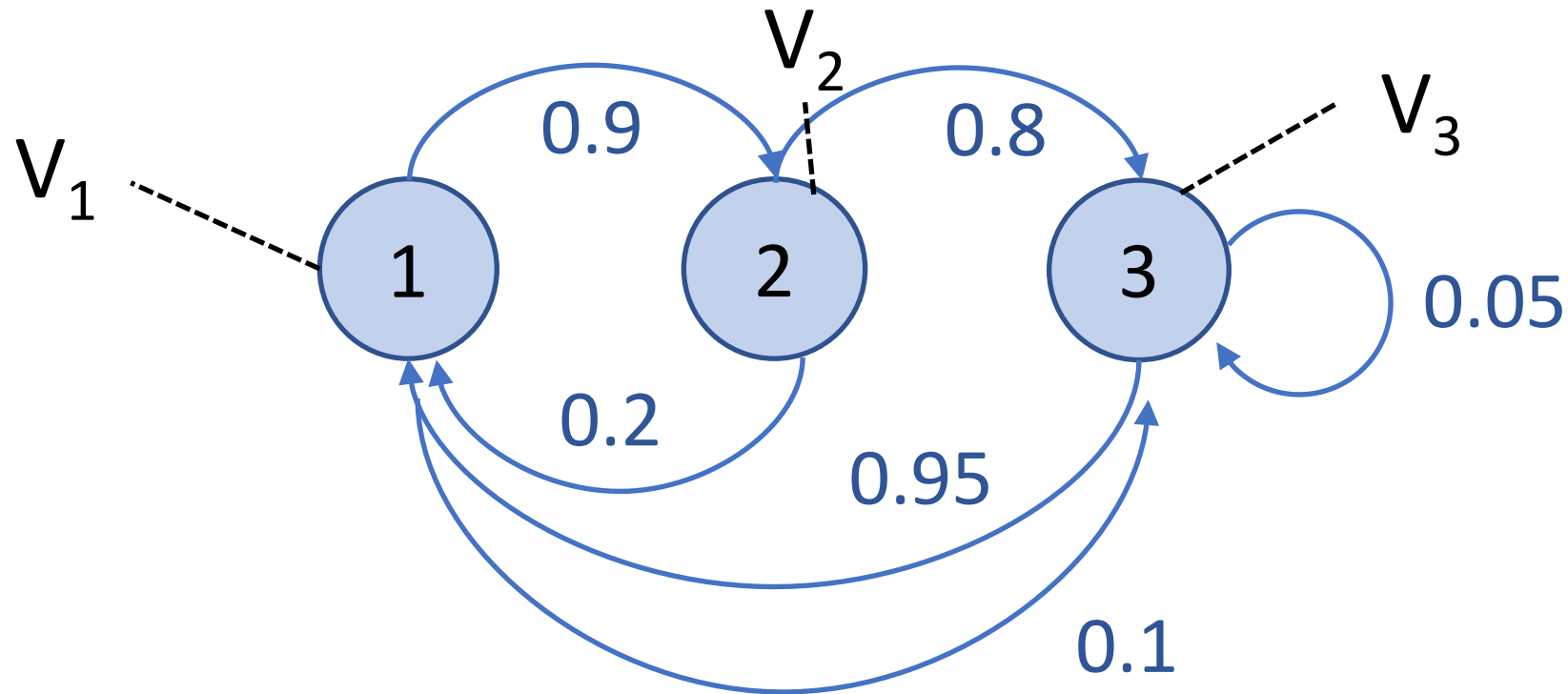
Hidden Markov Model

= Markov Network with Output Probabilities

# Markov Network (Chains)

= FSM with Transition Probabilities (a's)

States correspond to physical events,  
observations, measurements



**Markov property:** the future state of a system depends only on its current state, not on the sequence of events that preceded it

# Probability Axioms (Kolmogorov Axioms)

1. Probabilities are non-negative reals:

$$P(A) \geq 0 \quad \text{for all events } A$$

2. The entire event space has probability 1:

$$\sum P(A) = 1 \quad (\text{sum over all events } A)$$

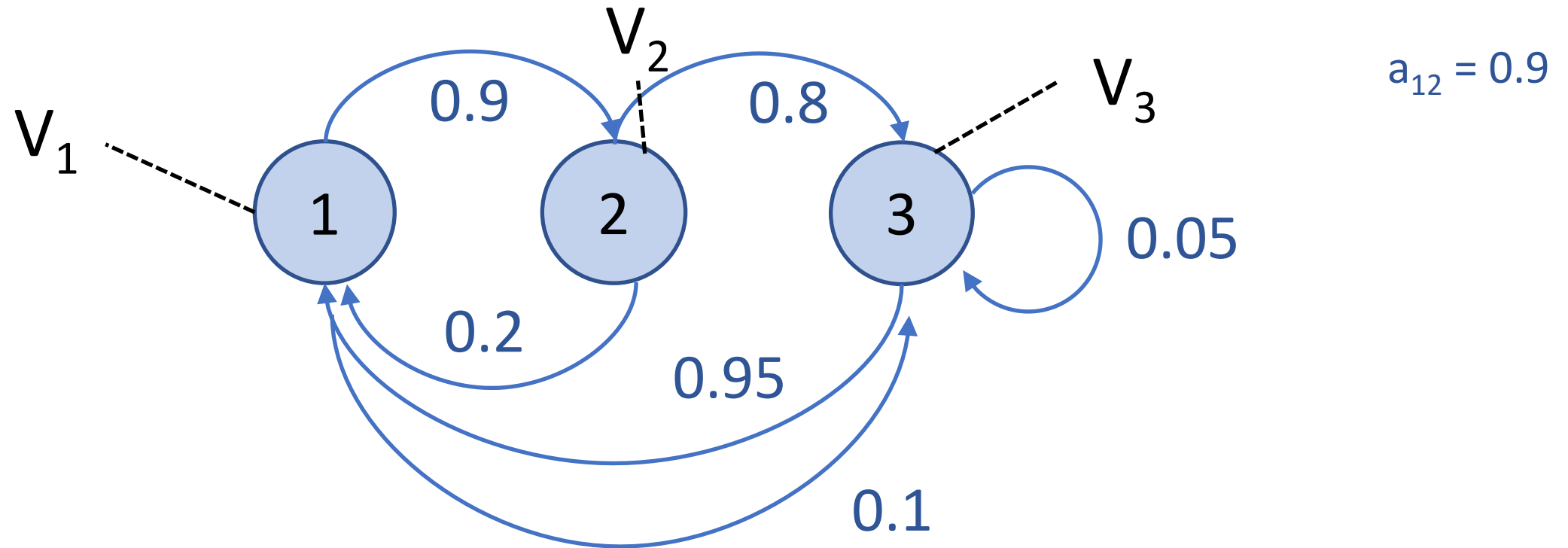
3. Probabilities add for pairwise disjoint events:

$$A \text{ \& B disjoint events: } P(A \cup B) = P(A) + P(B)$$

# Markov Chains

= FSM with Transition Probabilities ( $a$ 's)

States correspond to physical events,  
observations, measurements

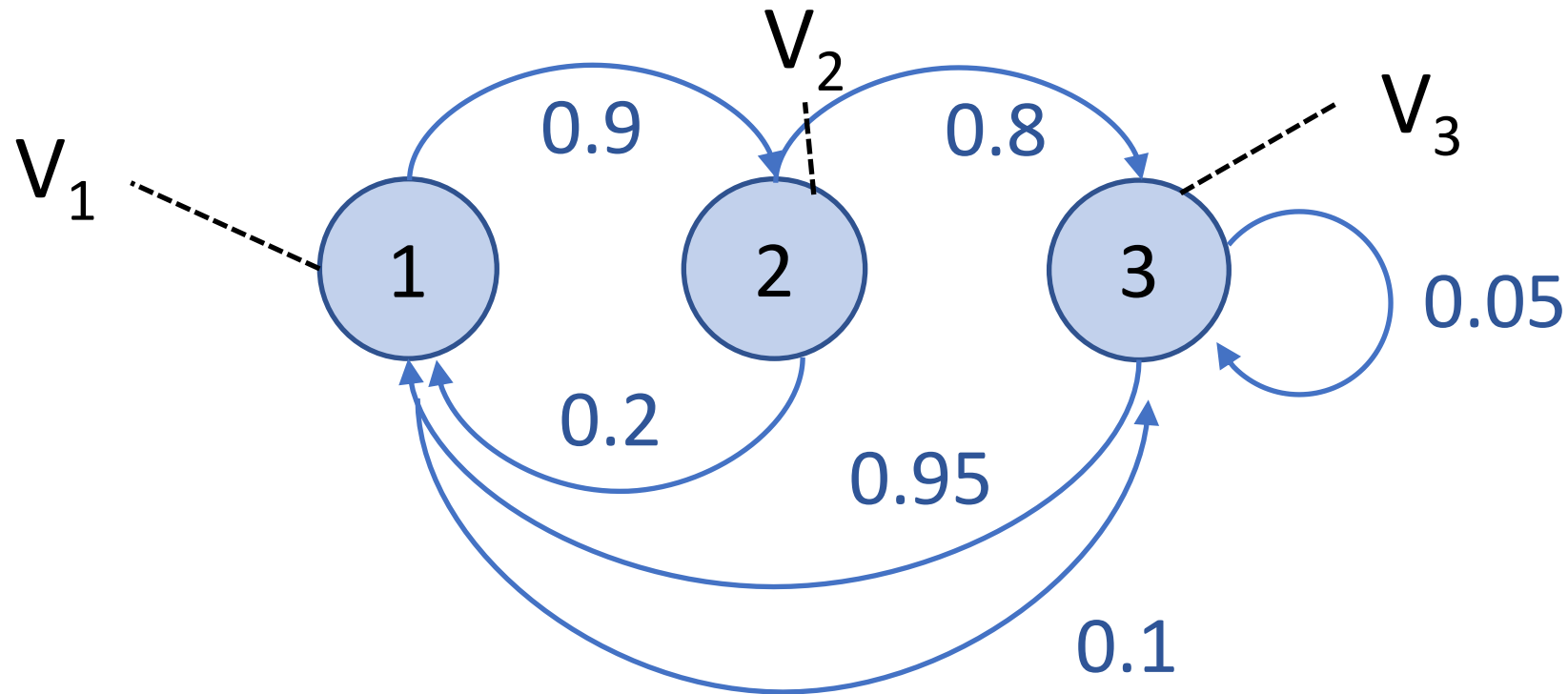


$$\sum_{j=1}^{\text{\#states}} \text{Prob( state } i \rightarrow \text{ state } j) = \sum_{j=1}^{\text{\#states}} a_{ij} = 1 \quad \text{for all states } i$$

# Markov Chains

= FSM with Transition Probabilities (a's)

States correspond to physical events,  
observations, measurements



$\pi_i$  = Probability that state  $i$  is the initial state

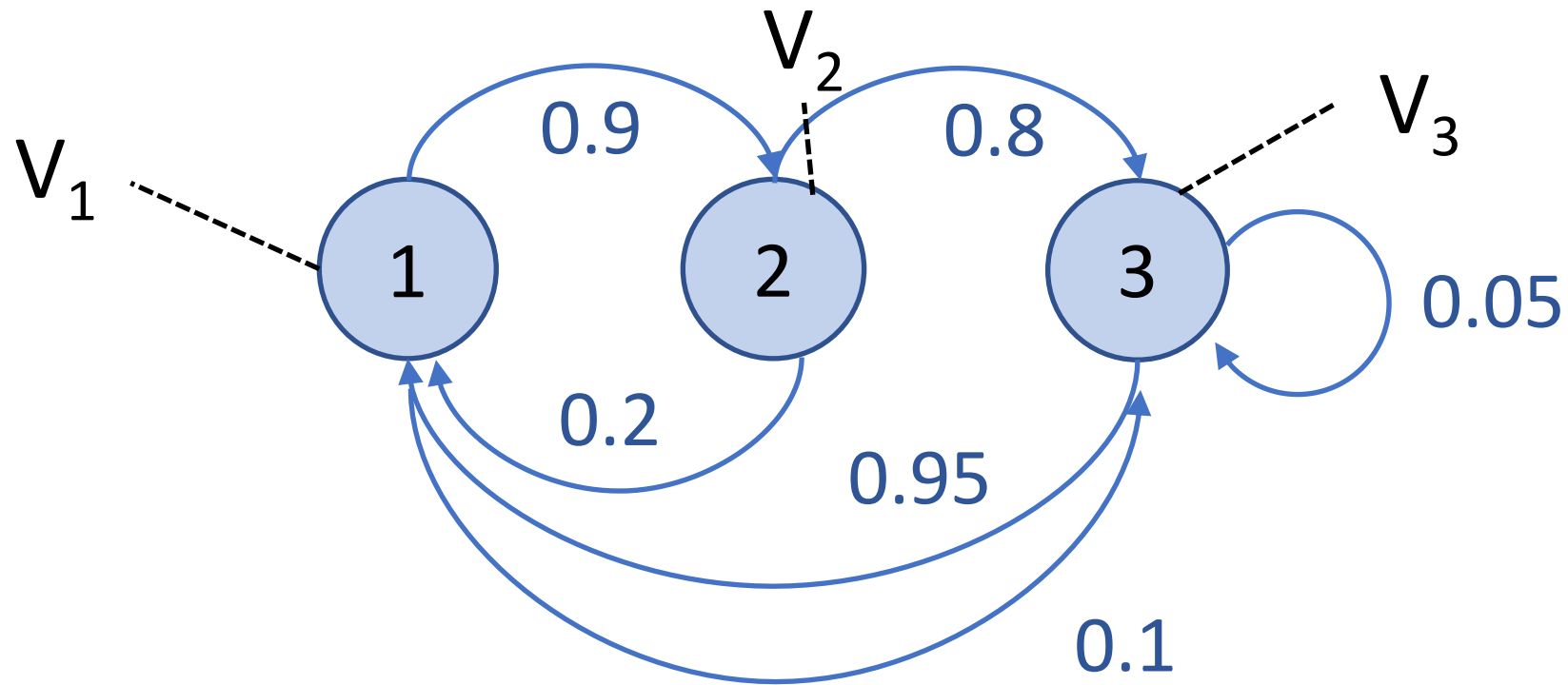
$$\pi_1 + \pi_2 + \pi_3 = 1$$

$$\pi_1 = 0.5$$

$$\pi_2 = 0.2$$

$$\pi_3 = 0.3$$

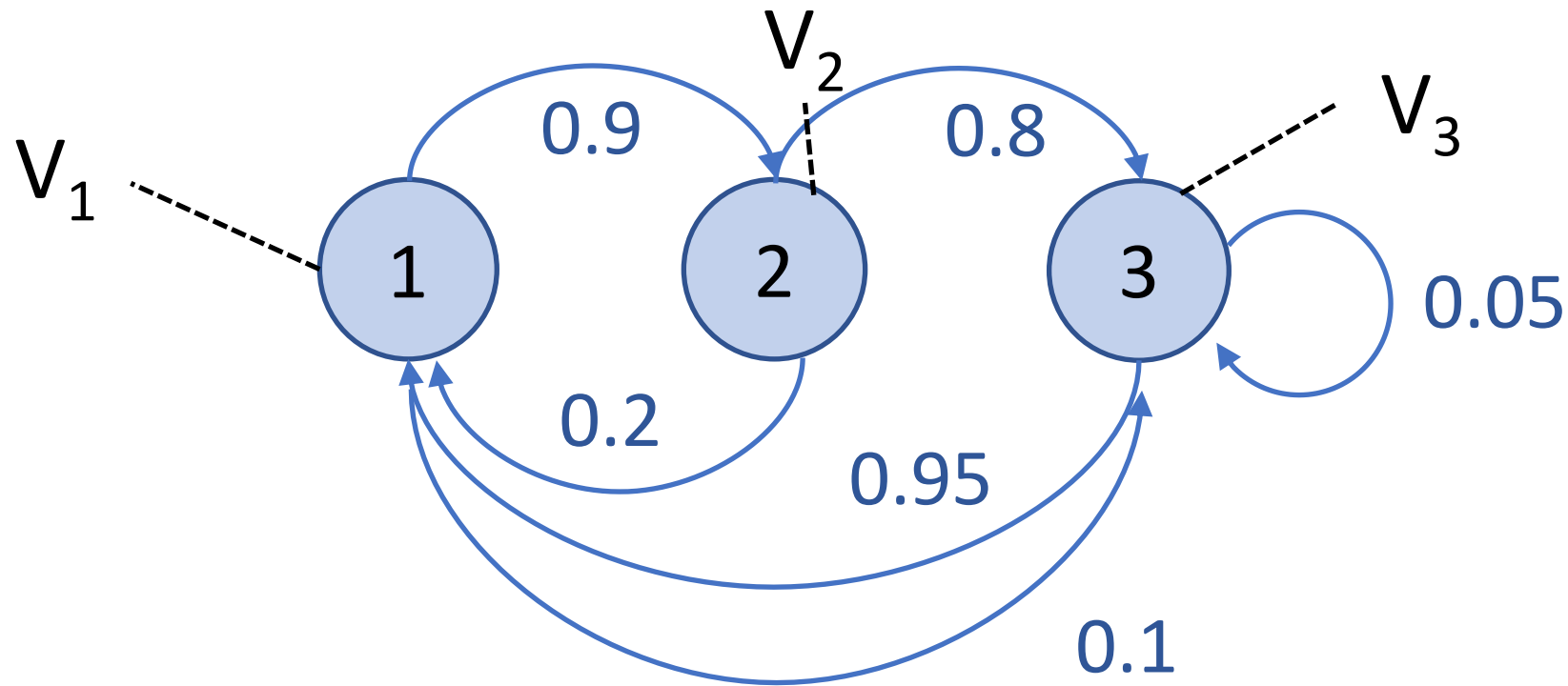
# Markov Chains



Given a sequence of outputs, we can calculate the probability of observing it.

E.g, outputs =  $V_2V_1V_3$

# Markov Chains

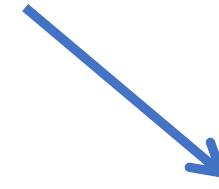


Probability of observing a sequence of outputs  $V_2V_1V_3 =$

$$\text{Prob}(O) = \text{Prob}(V_2V_1V_3) = \underset{0.2}{\pi_2} \cdot \underset{0.2}{\text{Prob}(\textcircled{2} \rightarrow \textcircled{1})} \cdot \underset{0.1}{\text{Prob}(\textcircled{1} \rightarrow \textcircled{3})} = 0.004$$

Finite State Machine (FSM)

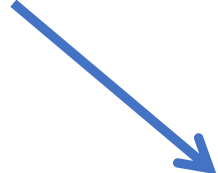
= Automata



Markov Network

= FSM with Transition Probabilities

Finite State Machine with  
Deterministic Outputs



Hidden Markov Model

= Markov Network with Output Probabilities



states not directly observable (= hidden)

# Hidden Markov Model

Markov model: each state outputs and must output one symbol, making the state outputs deterministic (**observable**).

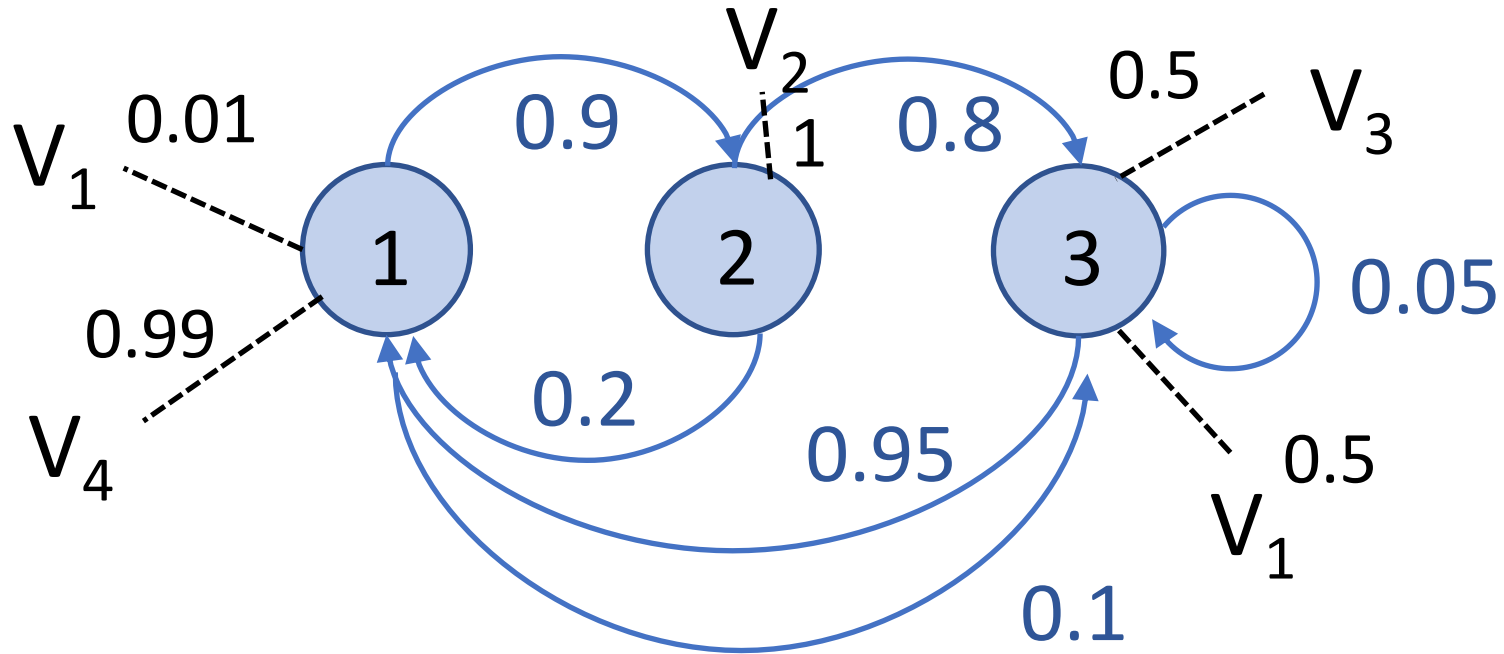
However, if instead, each state can output different symbols where each symbol is associated with a certain probability, then the outputs are non-deterministic (**hidden**) and the resulting model is called **Hidden Markov Model (HMM)**.

At a high level, a HMM is a Markov model with Markov transition process and non-observable (hidden) states.

# HMM Notations

- A set of  $N$  **states**:  $\mathbf{S} = \{S_1, S_2, S_3, \dots, S_N\}$ 
  - denote  $q_t$  : the actual state at time  $t$
- A set of  $M$  distinct symbols as **vocabulary** :  $\mathbf{V} = \{V_1, V_2, V_3, \dots, V_M\}$ 
  - sometimes symbols are in lower case
- Transition probabilities as a **matrix**:  $A = \{a_{ij}\}$ 
  - $a_{ij} = \text{Prob}(q_t = S_j \mid q_{t-1} = S_i), 1 \leq i, j \leq N$
- Initial probabilities:  $\boldsymbol{\pi} = \{\pi_1, \pi_2, \pi_3, \dots, \pi_N\}$ 
  - $\pi_i = \text{Prob}(q_1 = S_i), 1 \leq i \leq N$
- A **matrix** for observation likelihoods (aka. Emission probabilities):  $B = \{b_j(k)\}$ 
  - $b_j(k) = \text{Prob}(V_k \text{ at } t \mid q_t = S_j), 1 \leq j \leq N, 1 \leq k \leq M$
  - The probability of  $V_k$  being generated from a state  $q_j$
- A sequence of  $T$  observations (observed symbols at time  $T$ ):  $\mathbf{O} = O_1 O_2 O_3 \dots O_T$ 
  - Each  $O_i$  is drawn from the vocabulary  $\mathbf{V}$
- HMM:  $\boldsymbol{\lambda} = (\mathbf{S}, \mathbf{V}, A, B, \boldsymbol{\pi})$  or more commonly just  $(A, B, \boldsymbol{\pi})$

# HMM Example



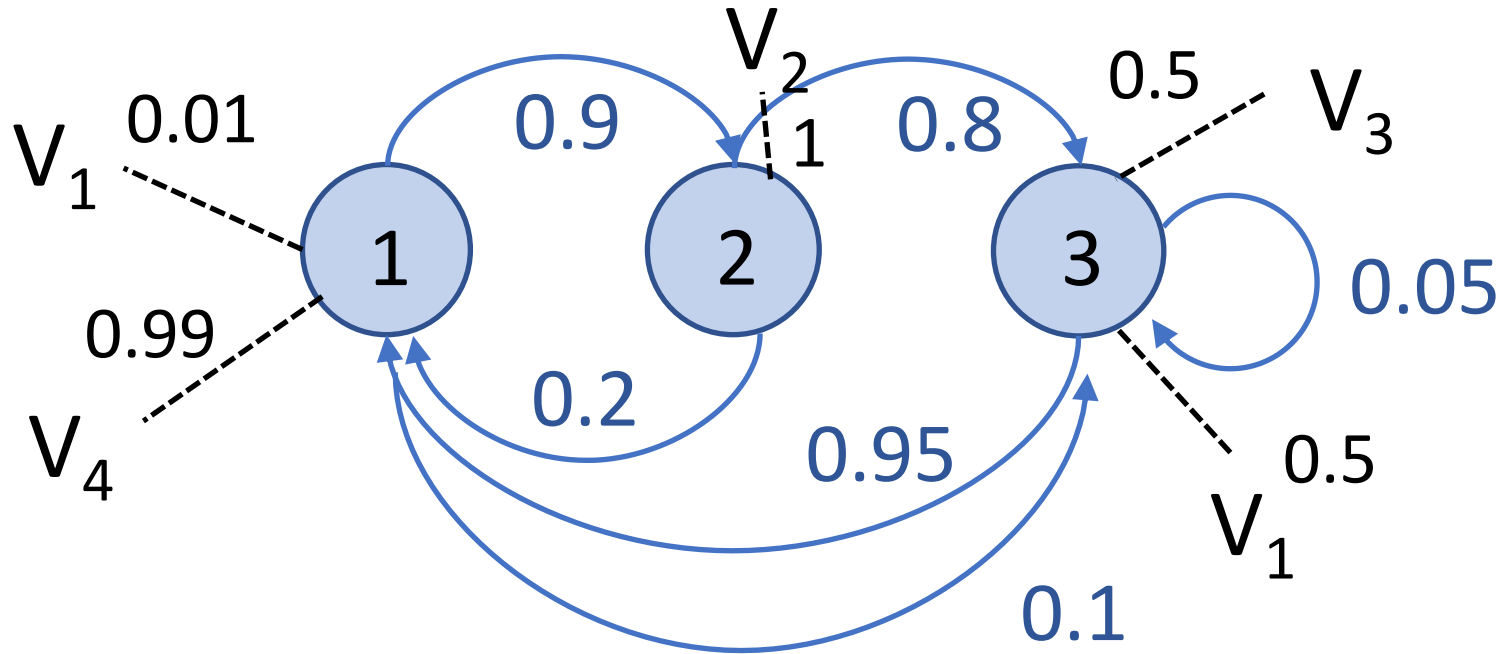
$S = ?$

$V = ?$

$A = ?$

$B = ?$

# HMM Example – States and Symbols



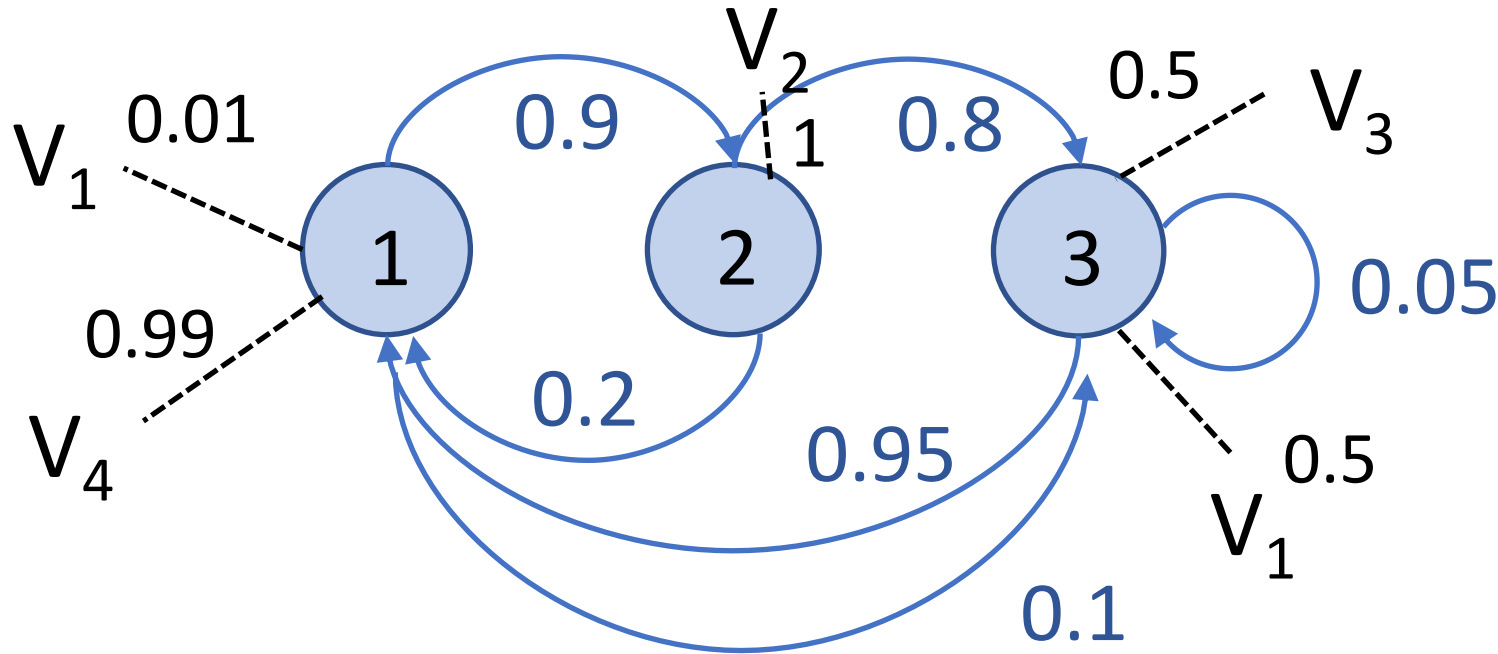
$S = [S_1, S_2, S_3]$

$V = [V_1, V_2, V_3, V_4]$

$A = ?$

$B = ?$

# HMM Example – Transition Probability



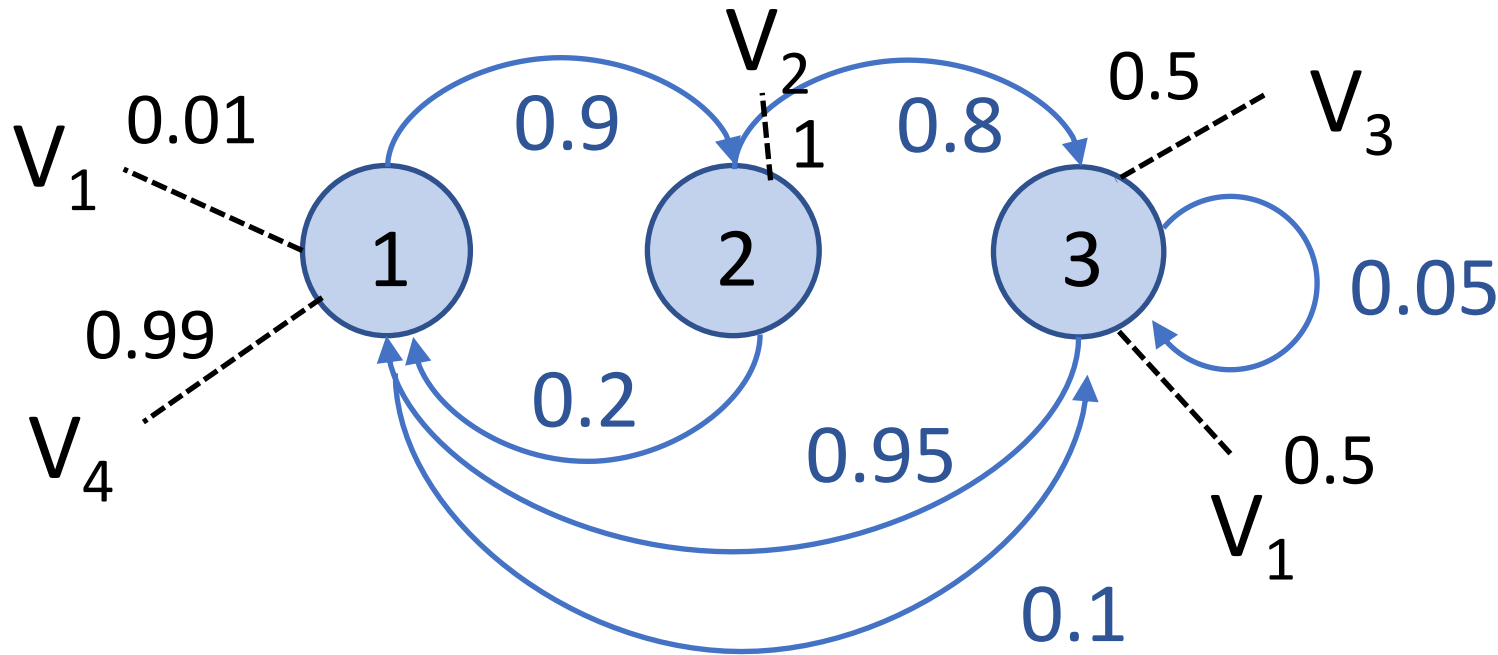
$S = [S_1, S_2, S_3]$

$V = [V_1, V_2, V_3, V_4]$

$A = \{a_{ij}\} = ?$

What is the size of  $A$ ?

# HMM Example – Transition Probability

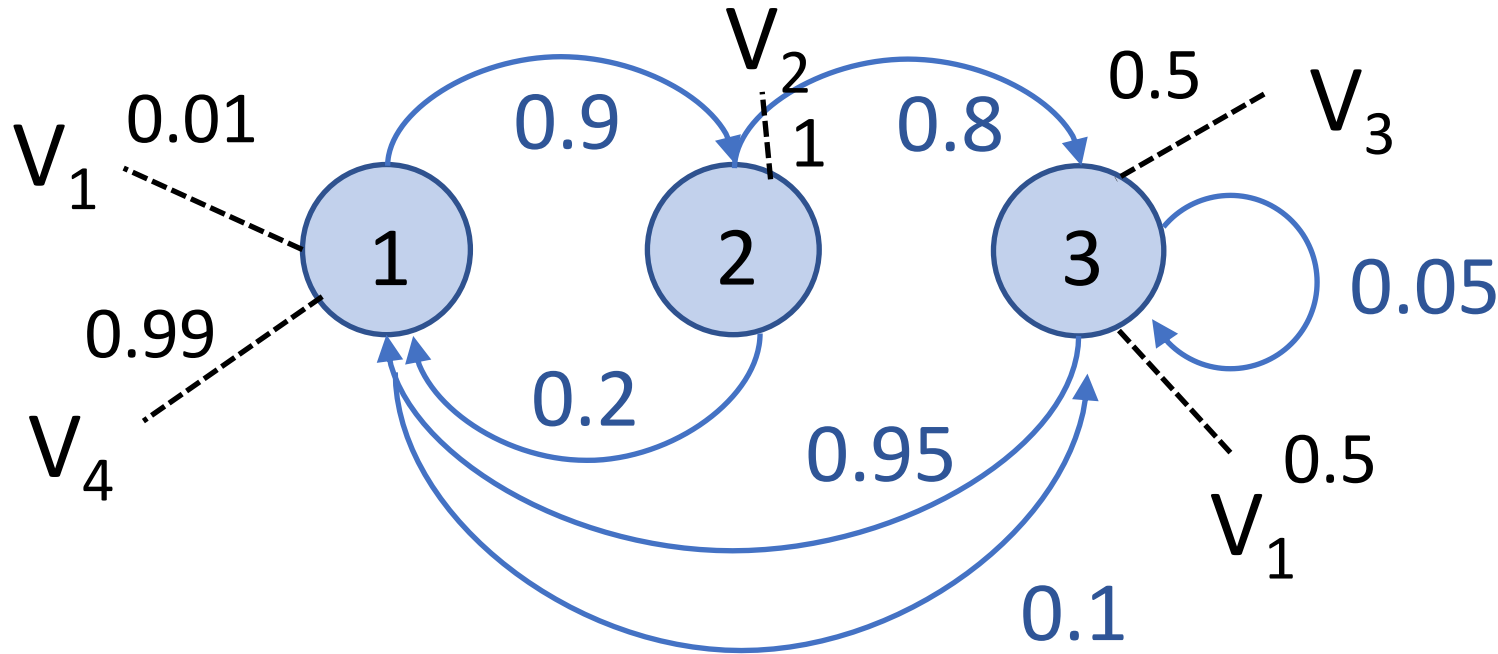


$$S = [S_1, S_2, S_3]$$

$$V = [V_1, V_2, V_3, V_4]$$

$$A = \{a_{ij}\}$$
$$= \begin{bmatrix} 0 & 0.9 & 0.1 \\ 0.2 & 0 & 0.8 \\ 0.95 & 0 & 0.05 \end{bmatrix}$$

# HMM Example – Emission Probability



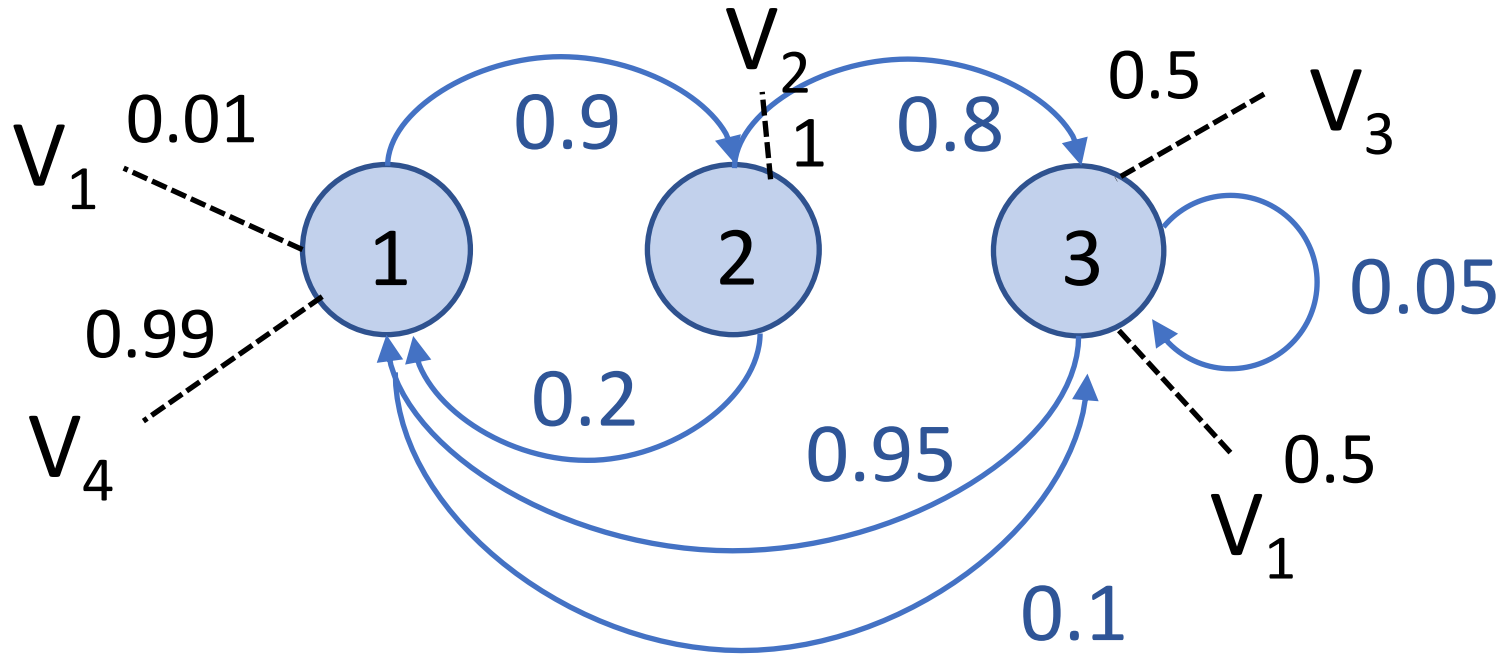
$$S = [S_1, S_2, S_3]$$

$$V = [V_1, V_2, V_3, V_4]$$

$$A = \{a_{ij}\}$$
$$= \begin{bmatrix} 0 & 0.9 & 0.1 \\ 0.2 & 0 & 0.8 \\ 0.95 & 0 & 0.05 \end{bmatrix}$$

$$B = ?$$

# HMM Example – Emission Probability



$$S = [S_1, S_2, S_3]$$

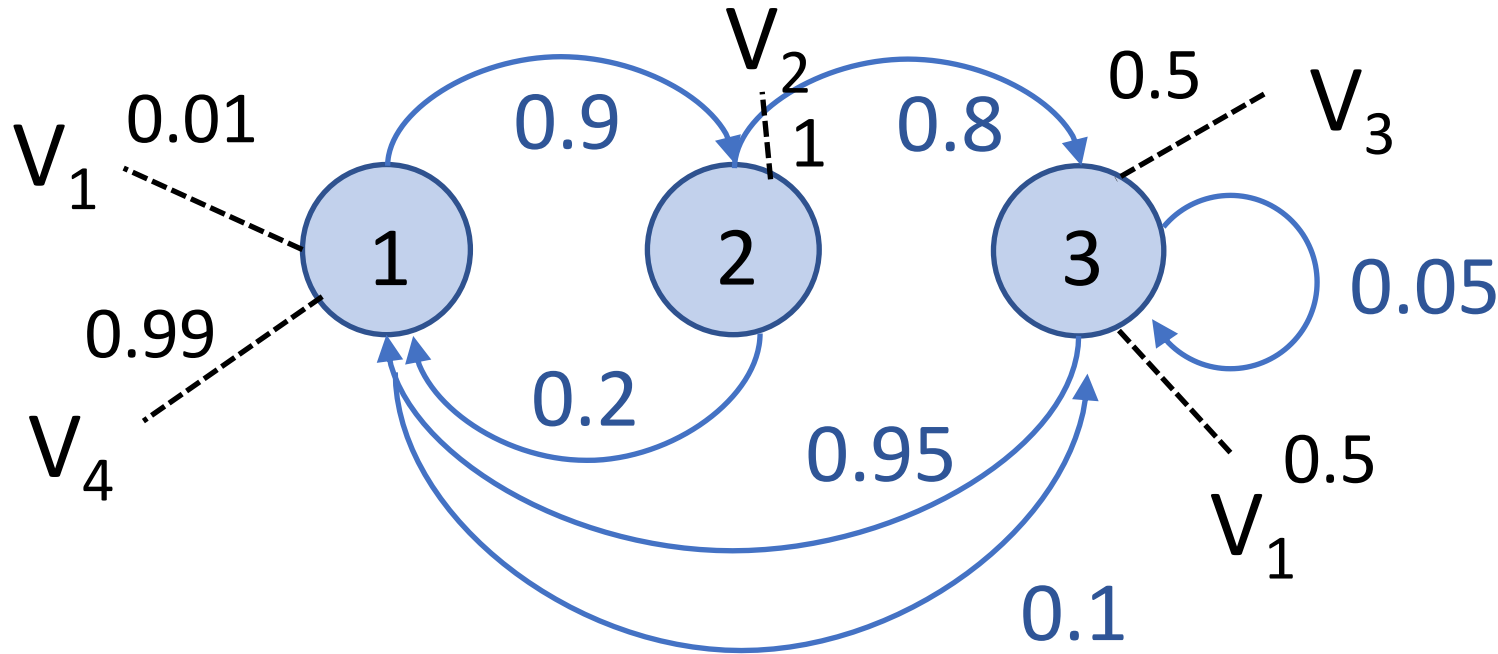
$$V = [V_1, V_2, V_3, V_4]$$

$$A = \{a_{ij}\}$$
$$= \begin{bmatrix} 0 & 0.9 & 0.1 \\ 0.2 & 0 & 0.8 \\ 0.95 & 0 & 0.05 \end{bmatrix}$$

$$B = [b_1, b_2, b_3] = ?$$



# HMM Example – Emission Probability



$$S = [S_1, S_2, S_3]$$

$$V = [V_1, V_2, V_3, V_4]$$

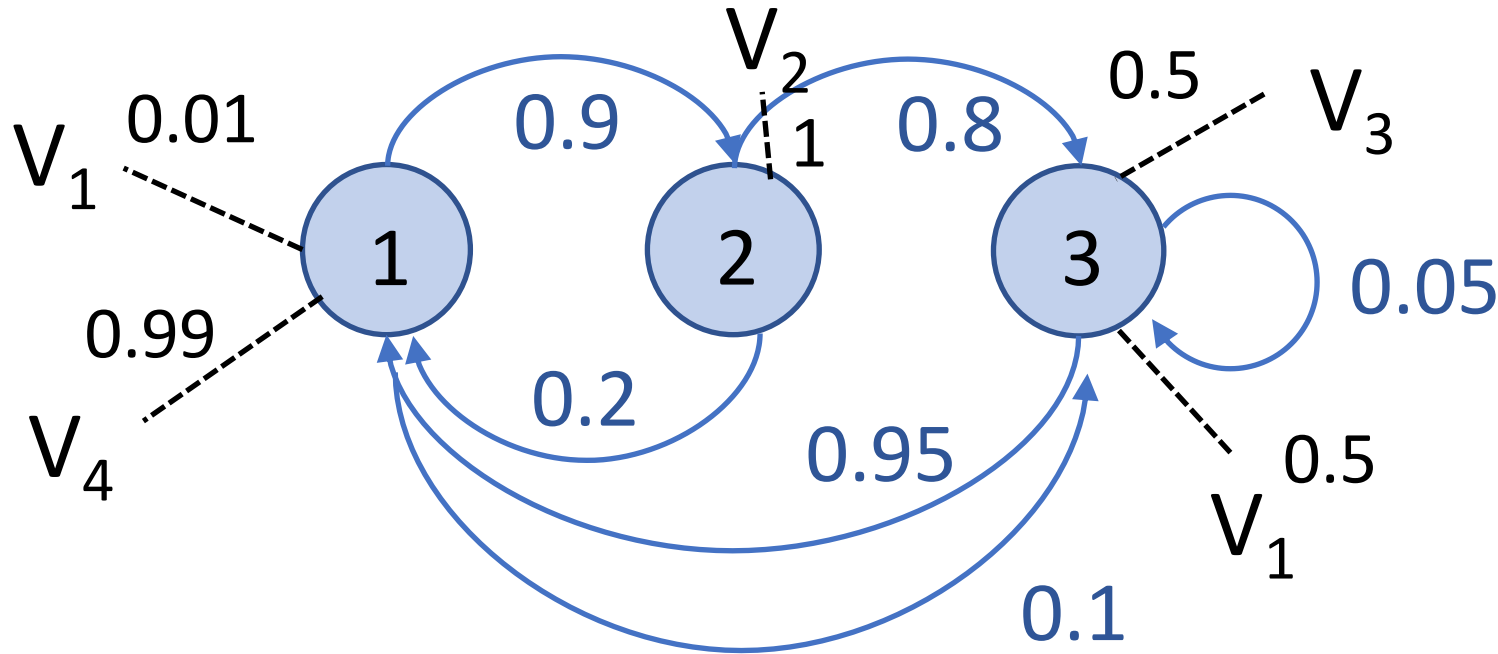
$$A = \{a_{ij}\}$$
$$= \begin{bmatrix} 0 & 0.9 & 0.1 \\ 0.2 & 0 & 0.8 \\ 0.95 & 0 & 0.05 \end{bmatrix}$$

$$B = [b_1, b_2, b_3]$$

where

$$b_i = [b_i(V_1), b_i(V_2), b_i(V_3), b_i(V_4)]$$

# HMM Example – Emission Probability



$$S = [S_1, S_2, S_3]$$

$$V = [V_1, V_2, V_3, V_4]$$

$$A = \{a_{ij}\}$$
$$= \begin{bmatrix} 0 & 0.9 & 0.1 \\ 0.2 & 0 & 0.8 \\ 0.95 & 0 & 0.05 \end{bmatrix}$$

$$B = [b_1, b_2, b_3]$$

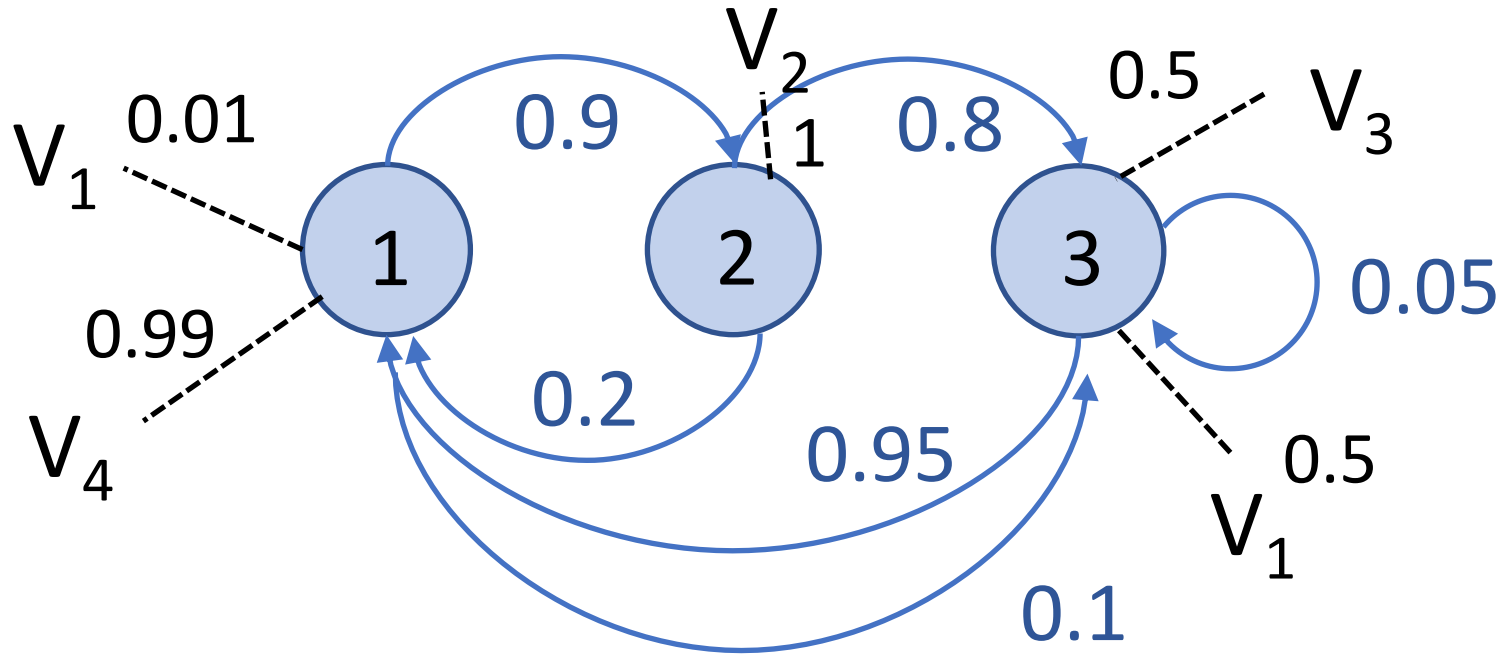
where

$$b_i = [b_i(V_1), b_i(V_2), b_i(V_3), b_i(V_4)]$$

where

$b_i(V_1)$  = the prob of observing  $V_1$  at state  $q_i$

# HMM Example – Emission Probability



$$Q = [q_1, q_2, q_3]$$

$$V = [V_1, V_2, V_3, V_4]$$

$$A = \{a_{ij}\}$$

$$= \begin{bmatrix} 0 & 0.9 & 0.1 \\ 0.2 & 0 & 0.8 \\ 0.95 & 0 & 0.05 \end{bmatrix}$$

$$B = [b_1, b_2, b_3]$$

where

$$b_i = [b_i(V_1), b_i(V_2), b_i(V_3), b_i(V_4)]$$

where

$b_i(V_1)$  = the prob of observing  $V_1$  at state  $q_i$

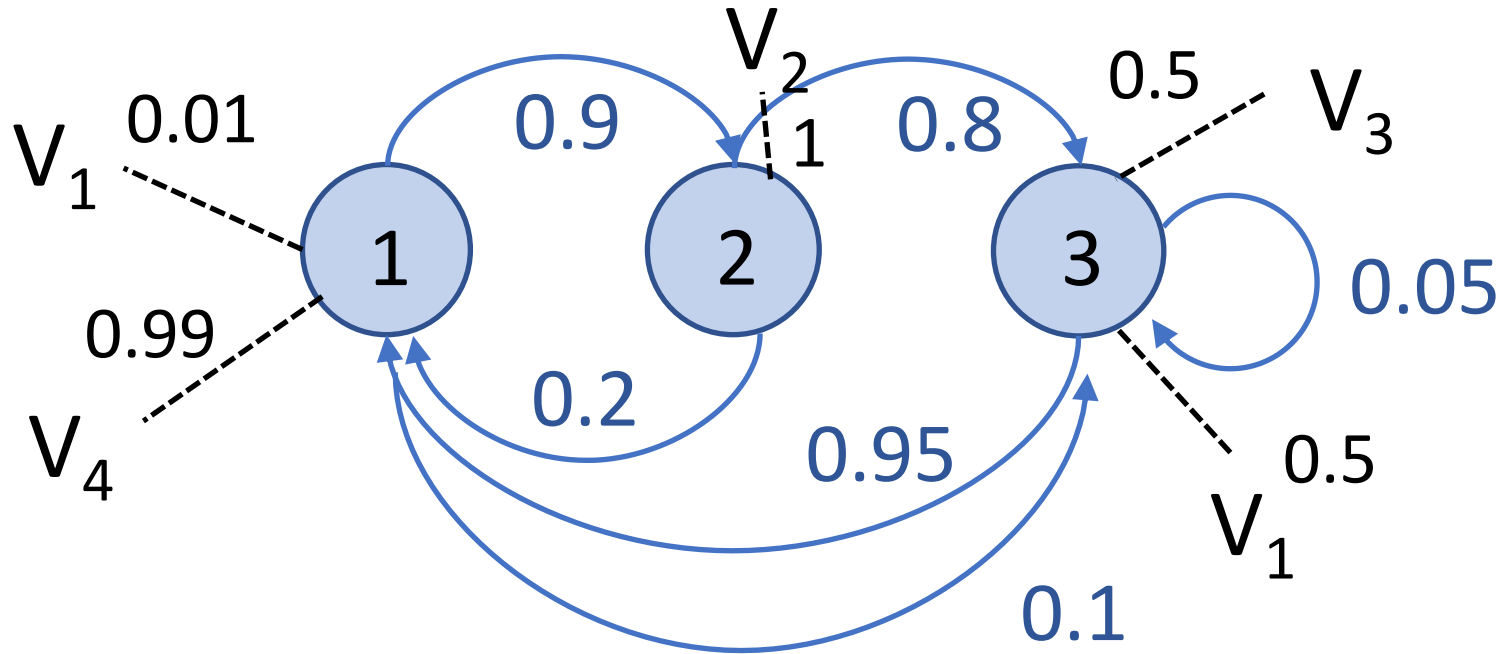
$$\sum_{k=1}^M \text{Prob (Output } V_k \text{ in state } i) = \sum_{k=1}^M b_i(V_k) = 1$$

$$b_1(V_1) = 0.01, b_1(V_2) = 0, b_1(V_3) = 0, b_1(V_4) = 0.99$$

$$b_2(V_1) = 0, b_2(V_2) = 1, b_2(V_3) = 0, b_2(V_4) = 0$$

$$b_3(V_1) = 0.5, b_3(V_2) = 0, b_3(V_3) = 0.5, b_3(V_4) = 0$$

# HMM Example – $\lambda$



$$S = [S_1, S_2, S_3]$$

$$V = [V_1, V_2, V_3, V_4]$$

$$A = \{a_{ij}\}$$
$$= \begin{bmatrix} 0 & 0.9 & 0.1 \\ 0.2 & 0 & 0.8 \\ 0.95 & 0 & 0.05 \end{bmatrix}$$

$$B = [b_1, b_2, b_3]$$

where

$$b_i = [b_i(V_1), b_i(V_2), b_i(V_3), b_i(V_4)]$$

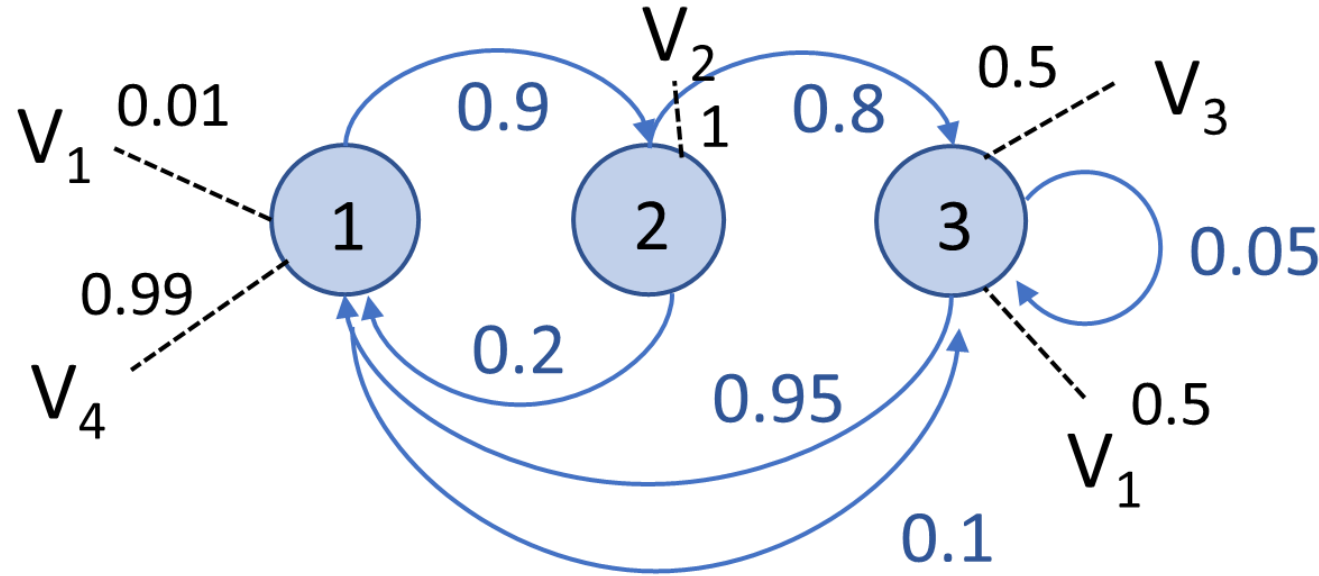
$$\lambda := (A, B, \pi)$$

Define specific HMM  $\lambda$ :

Transition matrix  $A$ , Emission probability  $B$ , Initial state vector  $\pi$

$$\pi = [0.5, 0.2, 0.3]$$

# HMM Example – Observation Probabilities

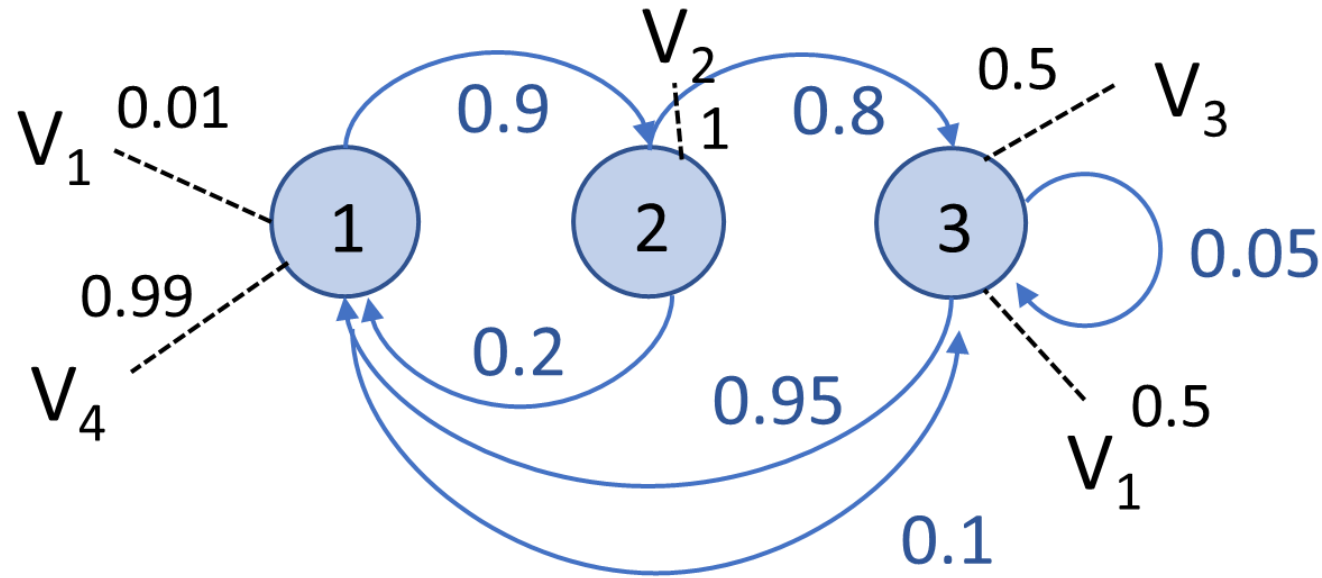


$$\begin{aligned}\pi_1 &= 0.5 \\ \pi_2 &= 0.2 \\ \pi_3 &= 0.3\end{aligned}$$

What is the probabilities of observing  $V_2$ ?

$$\begin{aligned}P(V_2) &= \pi_1 \cdot b_1(V_2) + \pi_2 \cdot b_2(V_2) + \pi_3 \cdot b_3(V_2) \\ &= 0 + 0.2 \cdot 1 + 0 \\ &= 0.2\end{aligned}$$

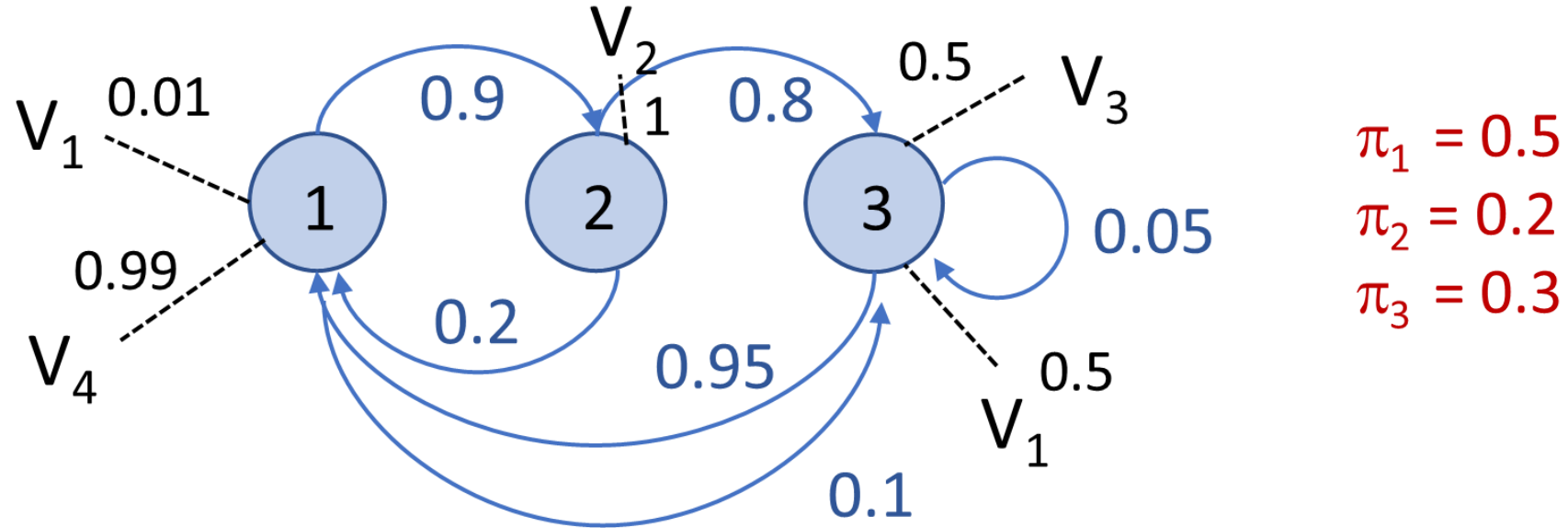
# HMM Example – Observation Probabilities



$$\begin{aligned}\pi_1 &= 0.5 \\ \pi_2 &= 0.2 \\ \pi_3 &= 0.3\end{aligned}$$

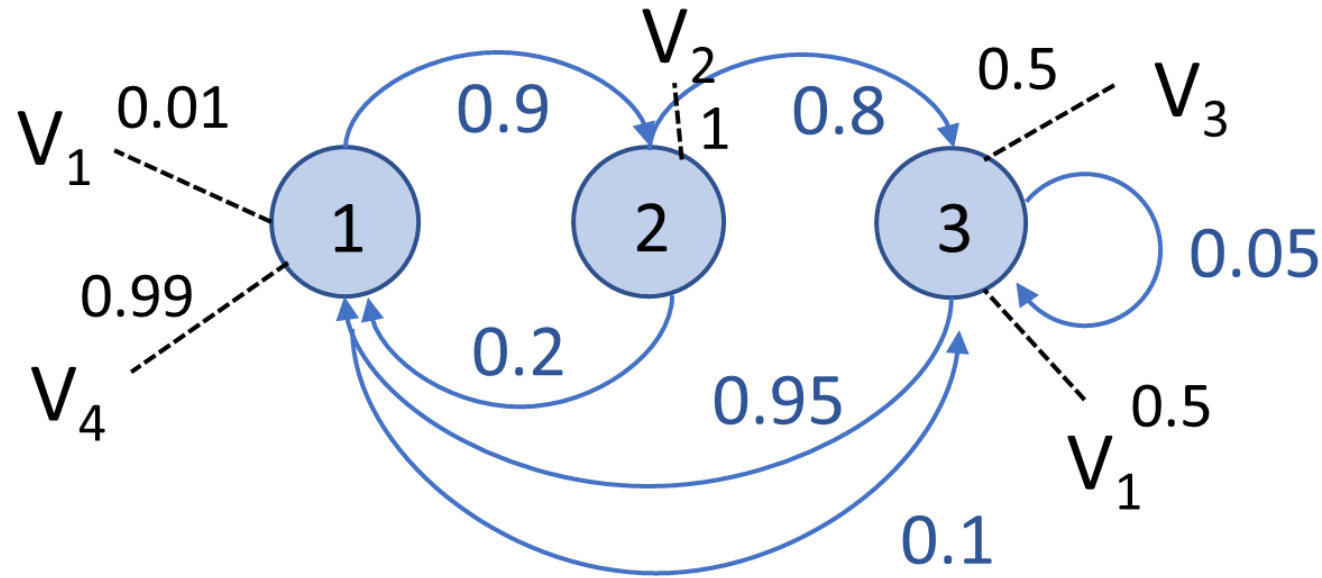
What is the probabilities of observing  $V_2 V_2$ ?

# HMM Example – Observation Probabilities



$$\begin{aligned} P(V_2 V_2) &= \pi_1 \cdot b_1(V_2) \cdot a_{11} \cdot b_1(V_2) + \pi_1 \cdot b_1(V_2) \cdot a_{12} \cdot b_2(V_2) + \pi_1 \cdot b_1(V_2) \cdot a_{13} \cdot b_3(V_2) + \\ &\quad \pi_2 \cdot b_2(V_2) \cdot a_{21} \cdot b_1(V_2) + \pi_2 \cdot b_2(V_2) \cdot a_{22} \cdot b_2(V_2) + \pi_2 \cdot b_2(V_2) \cdot a_{23} \cdot b_3(V_2) + \\ &\quad \pi_3 \cdot b_3(V_2) \cdot a_{31} \cdot b_1(V_2) + \pi_3 \cdot b_3(V_2) \cdot a_{32} \cdot b_2(V_2) + \pi_3 \cdot b_3(V_2) \cdot a_{33} \cdot b_3(V_2) \\ &= 0 + 0 + 0 \\ &= 0 \end{aligned}$$

# HMM Example – Observation Probabilities

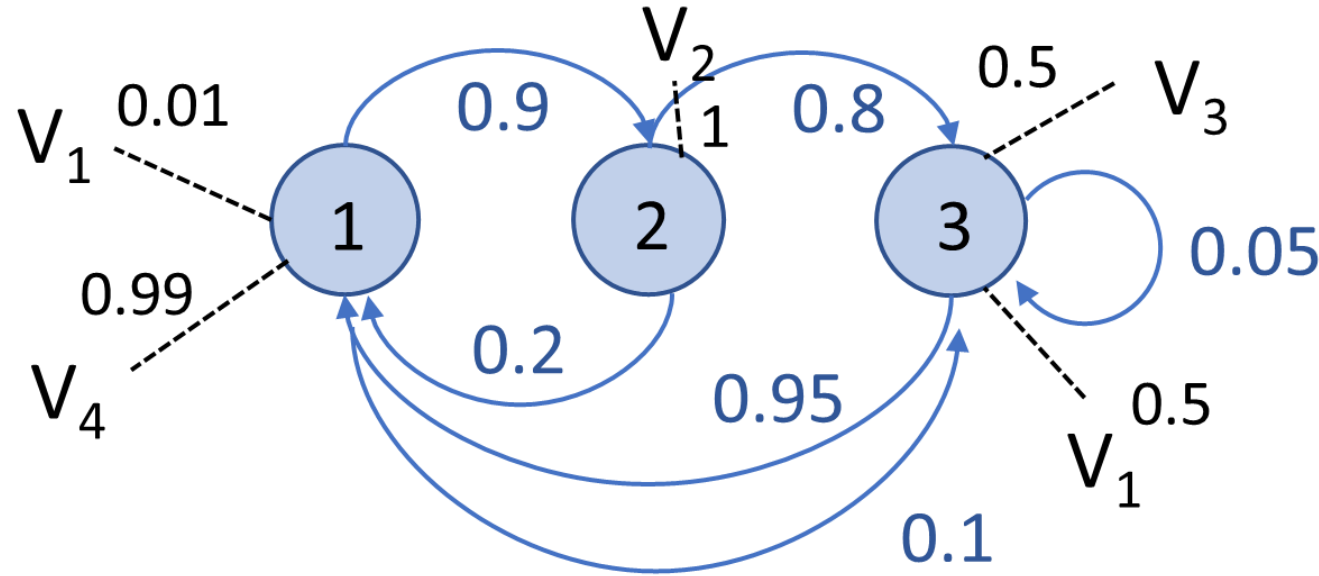


$$\begin{aligned}\pi_1 &= 0.5 \\ \pi_2 &= 0.2 \\ \pi_3 &= 0.3\end{aligned}$$

What is the probabilities of observing  $V_2 V_1$ ?



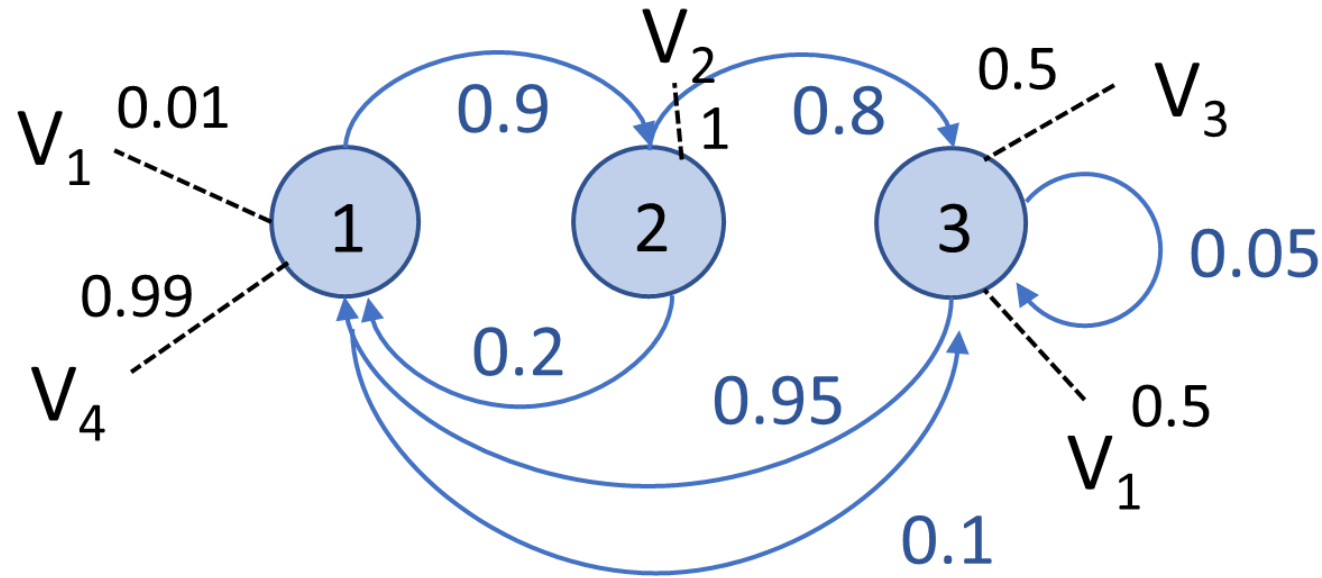
# HMM Example – Observation Probabilities



$$\begin{aligned}\pi_1 &= 0.5 \\ \pi_2 &= 0.2 \\ \pi_3 &= 0.3\end{aligned}$$

$$\begin{aligned}P(V_2V_1) &= \pi_1 \cdot b_1(V_2) \cdot a_{11} \cdot b_1(V_1) + \pi_1 \cdot b_1(V_2) \cdot a_{12} \cdot b_2(V_1) + \pi_1 \cdot b_1(V_2) \cdot a_{13} \cdot b_3(V_1) + \\ &\quad \pi_2 \cdot b_2(V_2) \cdot a_{21} \cdot b_1(V_1) + \pi_2 \cdot b_2(V_2) \cdot a_{22} \cdot b_2(V_1) + \pi_2 \cdot b_2(V_2) \cdot a_{23} \cdot b_3(V_1) + \\ &\quad \pi_3 \cdot b_3(V_2) \cdot a_{31} \cdot b_1(V_1) + \pi_3 \cdot b_3(V_2) \cdot a_{32} \cdot b_2(V_1) + \pi_3 \cdot b_3(V_2) \cdot a_{33} \cdot b_3(V_1) \\ &= 0 + (0.2 * 1 * 0.2 * 0.01 + 0 + 0.2 * 1 * 0.8 * 0.5) + 0 \\ &= 0.0804\end{aligned}$$

# HMM Example – Observation Probabilities



$$\begin{aligned}\pi_1 &= 0.5 \\ \pi_2 &= 0.2 \\ \pi_3 &= 0.3\end{aligned}$$

What is the probabilities of observing  $V_2 V_1 V_3$ ?

# Working with Hidden Markov Models

There are three fundamental problems in HMM:

- **Evaluation Problem:**

How likely is it that HMM  $\lambda$  computed  $O$ ?

Forward or backward procedure

- **Recognition Problem:**

Does HMM  $\lambda$  recognize  $O$ ?

Viterbi Algorithm

- **Learning (= Training) Problem:**

Adjust  $\lambda$  so that  $\text{Prob}(O|\lambda)$  is locally maximized.

# Problem 1: Evaluation Problem

Problem Definition: Given an observation sequence  $O=O_1O_2O_3...O_T$  and the model  $\lambda$ , what is  $P(O|\lambda)$ ?

- The most straightforward way: enumerate every possible state sequence of length  $T$ 
  - This is what we have just practiced and we have seen how fast the math gets ugly.

# Problem 1: Evaluation Problem

Problem Definition: Given an observation sequence  $O=O_1O_2...O_T$  and the model  $\lambda$ , what is  $P(O|\lambda)$ ?

- Consider **one** such fixed state sequence  $Q=q_1q_2...q_T$ , we have
  - $P(O|Q, \lambda) = \prod_{t=1}^T P(O_t|q_t, \lambda) = b_{q_1}(O_1) \cdot b_{q_2}(O_2) \cdot \dots \cdot b_{q_T}(O_T)$
- The prob. of such a state seq  $Q$  can be written as
  - $P(Q|\lambda) = \pi_{q_1} a_{q_1q_2} a_{q_2q_3} \dots a_{q_{T-1}q_T}$
- The joint prob. of  $O$  and  $Q$ :
  - $P(O, Q|\lambda) = P(O|Q, \lambda)P(Q|\lambda)$
- With **all** possible state sequences, summing up
  - $P(O|\lambda) = \sum_{all\ Q} P(O, Q|\lambda) = \sum_{all\ Q} P(O|Q, \lambda)P(Q|\lambda)$

# Problem 1: Evaluation Problem

Problem Definition: Given an observation sequence  $O=O_1O_2O_3...O_T$  and the model  $\lambda$ , what is  $P(O|\lambda)$ ?

- The most straightforward way: enumerate every possible state sequence of length  $T$ 
  - This is what we have just practiced and we have seen how fast the math gets ugly.
  - The time complexity is exponential.  $O(2^T N^T)$
- More efficient procedure?

## Forward Procedure

- Define forward variable  $\alpha_t(i)$ :

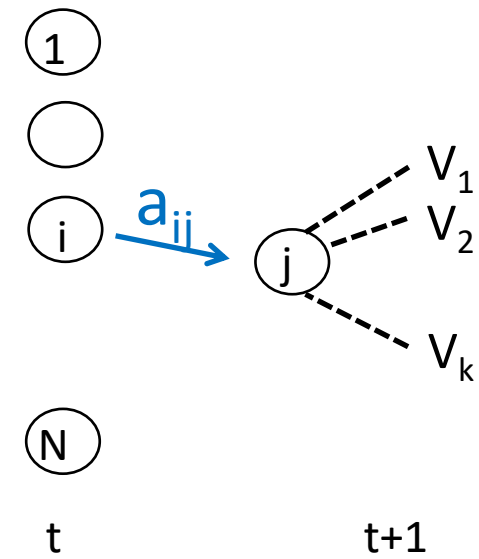
$$\alpha_t(i) = P(O_1 O_2 \dots O_t, q_t = S_i | \lambda)$$

$$O = O_1 O_2 \dots O_{t-1} O_t O_{t+1} \dots O_T$$

- Given the model,  $\alpha_t(i)$  is the probability of the **partial** observation sequence  $O_1 O_2 \dots O_t$  when it reaches state  $S_i$  at time  $t$ .
- How about  $\alpha_{t+1}(j)$  for some state  $S_j$ ?
  - assuming  $\alpha_t(i)$  is known

$$\alpha_{t+1}(j) = \left( \sum_i \alpha_t(i) \cdot a_{ij} \right) \cdot b_j(O_{t+1})$$

- We can solve  $\alpha_t(i)$  inductively.
- Time complexity.  $O(T N^2)$



## Forward Procedure

### 1. Initialization

$$\alpha_1(i) = \pi_i b_i(o_1), \quad 1 \leq i \leq N$$

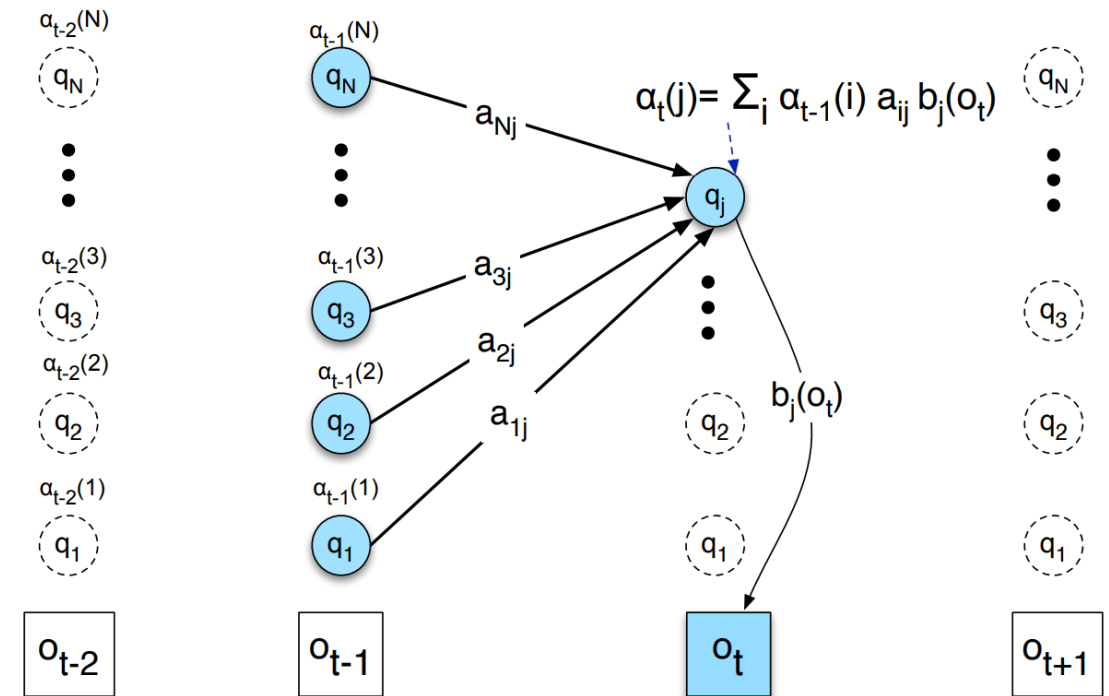
### 2. Induction

For  $t = 1, 2, \dots, T-1$

$$\alpha_{t+1}(j) = \left[ \sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(o_{t+1}), \quad 1 \leq j \leq N$$

### 3. Termination

$$\text{Prob}(O \mid \lambda) = \sum_{i=1}^N \alpha_T(i)$$



$$\alpha_t(i) = P(o_1 o_2 \dots o_t \& i \mid \lambda)$$



## Forward Procedure

### 1. Initialization

$$\alpha_1(i) = \pi_i b_i(o_1), \quad 1 \leq i \leq N$$

### 2. Induction

For  $t = 1, 2, \dots, T-1$

$$\alpha_{t+1}(j) = \left[ \sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(o_{t+1}), \quad 1 \leq j \leq N$$

### 3. Termination

$$\text{Prob}(O \mid \lambda) = \sum_{i=1}^N \alpha_T(i)$$

$\alpha$	1	2	3						T
1									
2									
3									
4									
N									

$$\alpha_t(i) = P(o_1 o_2 \dots o_t \& i \mid \lambda)$$

## Example: Forward Procedure

### 1. Initialization

$$\alpha_1(i) = \pi_i b_i(o_1), \quad 1 \leq i \leq N$$

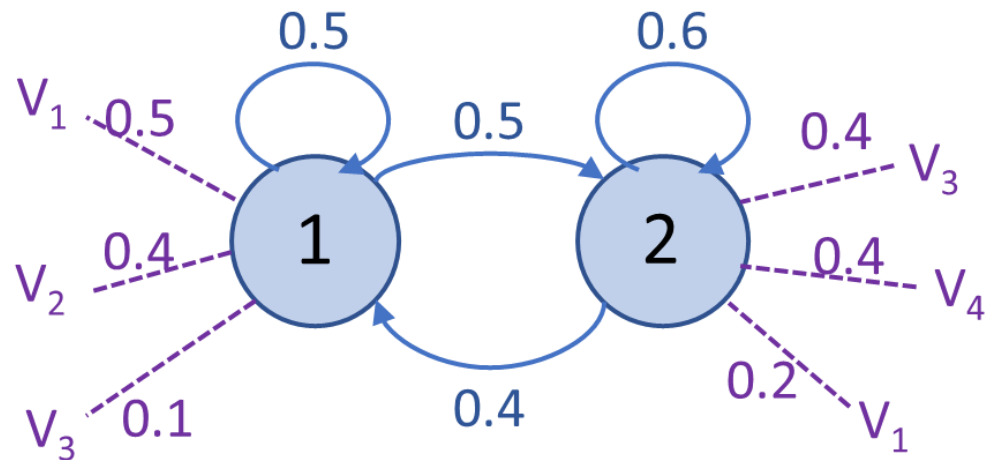
### 2. Induction

For  $t = 1, 2, \dots, T-1$

$$\alpha_{t+1}(j) = \left[ \sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(o_{t+1}), \quad 1 \leq j \leq N$$

### 3. Termination

$$\text{Prob}(O \mid \lambda) = \sum_{i=1}^N \alpha_T(i)$$



$$\pi_1 = 0.2, \pi_2 = 0.8$$

$\alpha$	1	2	3
1			
2			

Observation  $O = V_3 V_1 V_4$  – What is the probability?

## Example: Forward Procedure

### 1. Initialization

$$\alpha_1(i) = \pi_i b_i(o_1), \quad 1 \leq i \leq N$$

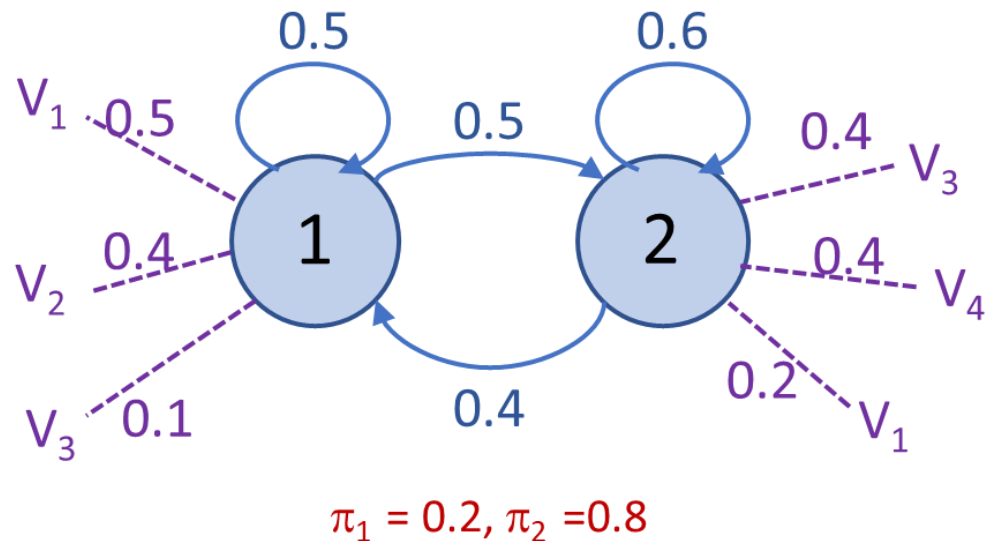
### 2. Induction

For  $t = 1, 2, \dots, T-1$

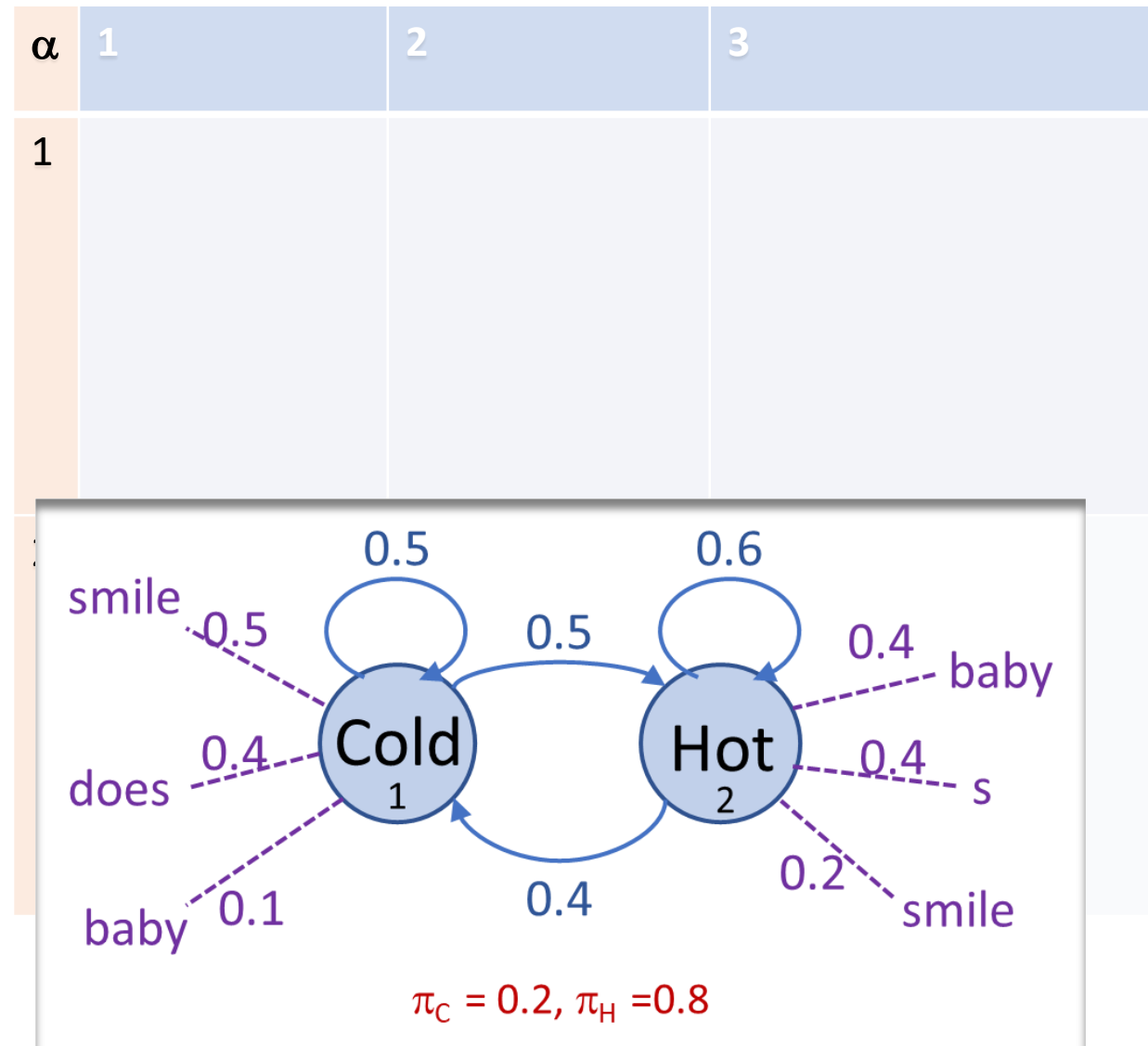
$$\alpha_{t+1}(j) = \left[ \sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(o_{t+1}), \quad 1 \leq j \leq N$$

### 3. Termination

$$\text{Prob}(O \mid \lambda) = \sum_{i=1}^N \alpha_T(i)$$



Observation  $O = V_3 V_1 V_4$  – What is the probability?



Observation  $O = V_3 V_1 V_4 \rightarrow O = \text{"baby smile s"}$

## Example: Forward Procedure

### 1. Initialization

$$\alpha_1(i) = \pi_i b_i(o_1), \quad 1 \leq i \leq N$$

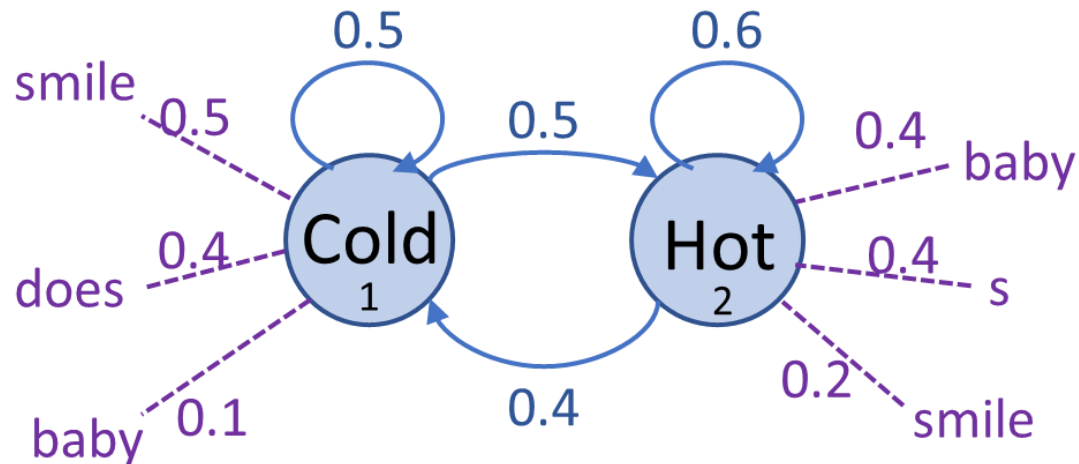
### 2. Induction

For  $t = 1, 2, \dots, T-1$

$$\alpha_{t+1}(j) = \left[ \sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(o_{t+1}), \quad 1 \leq j \leq N$$

### 3. Termination

$$\text{Prob}(O | \lambda) = \sum_{i=1}^N \alpha_T(i)$$



$$\pi_C = 0.2, \pi_H = 0.8$$

Observation  $O = \text{"baby smile s"} - \text{What is the probability?}$

$\alpha$	1	2	3
1 Cold	$\alpha_1(C)$ $= \pi_C b_C(\text{baby})$ $= 0.2 * 0.1$ $= 0.02$	$\alpha_2(C)$ $= (\alpha_1(C) \cdot a_{CC}$ $+ \alpha_1(H) \cdot a_{HC}) \cdot b_C(\text{smile})$ $= (0.02 * 0.5 + 0.32$ $* 0.4) * 0.5 = 0.069$	$\alpha_3(C)$ $= (\alpha_2(C) \cdot a_{CC}$ $+ \alpha_2(H) \cdot a_{HC}) \cdot b_C(s)$ $= (0.069 * 0.5$ $+ 0.0404 * 0.4) * 0$ $= 0$
2 Hot	$\alpha_1(H)$ $= \pi_H b_H(\text{baby})$ $= 0.8 * 0.4$ $= 0.32$	$\alpha_2(H)$ $= (\alpha_1(C) \cdot a_{CH}$ $+ \alpha_1(H) \cdot a_{HH}) \cdot b_H(\text{smile})$ $= (0.02 * 0.5 + 0.32$ $* 0.6) * 0.2 = 0.0404$	$\alpha_3(H)$ $= (\alpha_2(C) \cdot a_{CH}$ $+ \alpha_2(H) \cdot a_{HH}) \cdot b_H(s)$ $= (0.069 * 0.5$ $+ 0.0404 * 0.6) * 0.4$ $= 0.023496$

$$P(O|\lambda) = \alpha_3(C) + \alpha_3(H) = 0.023496$$

## Example: Forward Procedure

### 1. Initialization

$$\alpha_1(i) = \pi_i b_i(o_1), \quad 1 \leq i \leq N$$

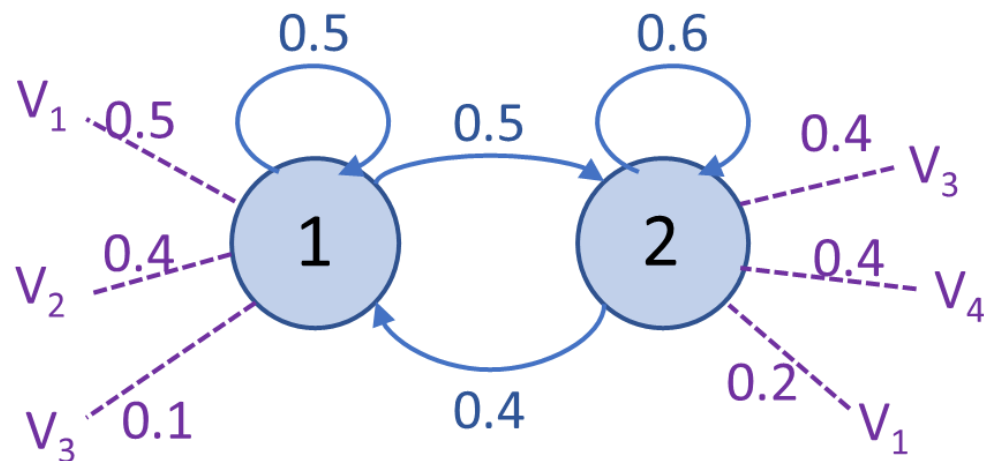
### 2. Induction

For  $t = 1, 2, \dots, T-1$

$$\alpha_{t+1}(j) = \left[ \sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(o_{t+1}), \quad 1 \leq j \leq N$$

### 3. Termination

$$\text{Prob}(O | \lambda) = \sum_{i=1}^N \alpha_T(i)$$



$$\pi_1 = 0.2, \pi_2 = 0.8$$

$\alpha$	1	2	3
1	$\begin{aligned} \alpha_1(1) &= \pi_1 b_1(V_3) \\ &= 0.2 * 0.1 \\ &= 0.02 \end{aligned}$	$\begin{aligned} \alpha_2(1) &= (\alpha_1(1) \cdot a_{11} + \alpha_1(2) \\ &\quad \cdot a_{21}) \cdot b_1(V_1) \end{aligned}$	$\begin{aligned} \alpha_3(1) &= (\alpha_2(1) \cdot a_{11} + \alpha_2(2) \\ &\quad \cdot a_{21}) \cdot b_1(V_4) \end{aligned}$
2	$\begin{aligned} \alpha_1(2) &= \pi_2 b_2(V_3) \\ &= 0.8 * 0.4 \\ &= 0.32 \end{aligned}$	$\begin{aligned} \alpha_2(2) &= (\alpha_1(1) \cdot a_{12} + \alpha_1(2) \\ &\quad \cdot a_{22}) \cdot b_2(V_1) \end{aligned}$	$\begin{aligned} \alpha_3(2) &= (\alpha_2(1) \cdot a_{12} + \alpha_2(2) \\ &\quad \cdot a_{22}) \cdot b_2(V_4) \end{aligned}$

$$P(O|\lambda) = \alpha_3(1) + \alpha_3(2)$$

Observation  $O = V_3 V_1 V_4$  – What is the probability?

## Example: Forward Procedure

### 1. Initialization

$$\alpha_1(i) = \pi_i b_i(o_1), \quad 1 \leq i \leq N$$

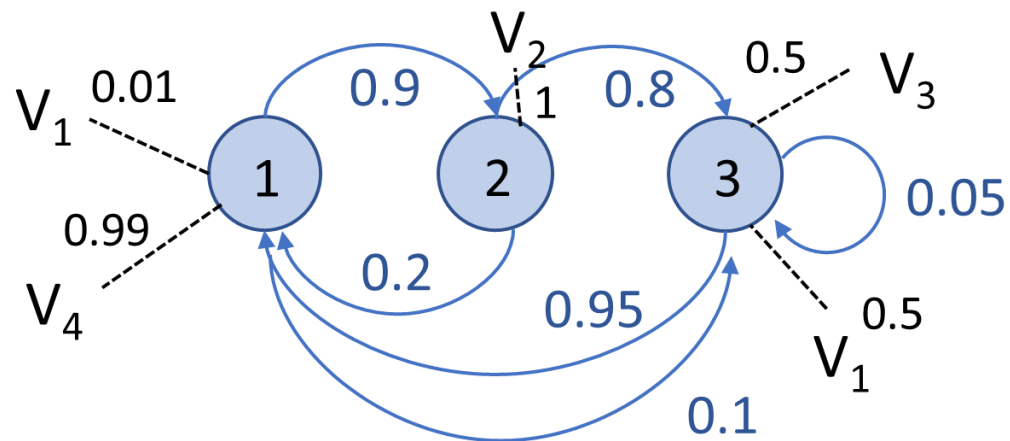
### 2. Induction

For  $t = 1, 2, \dots, T-1$

$$\alpha_{t+1}(j) = \left[ \sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(o_{t+1}), \quad 1 \leq j \leq N$$

### 3. Termination

$$\text{Prob}(O \mid \lambda) = \sum_{i=1}^N \alpha_T(i)$$



$$\pi_1 = 0.5$$

$$\pi_2 = 0.2$$

$$\pi_3 = 0.3$$

$$P(O=V_2V_1V_3) = ?$$

$\alpha$	1	2	3
1	$\alpha_1(1)$	$\alpha_2(1)$	$\alpha_3(1)$
2	$\alpha_1(2)$	$\alpha_2(2)$	$\alpha_3(2)$
3	$\alpha_1(3)$	$\alpha_2(3)$	$\alpha_3(3)$

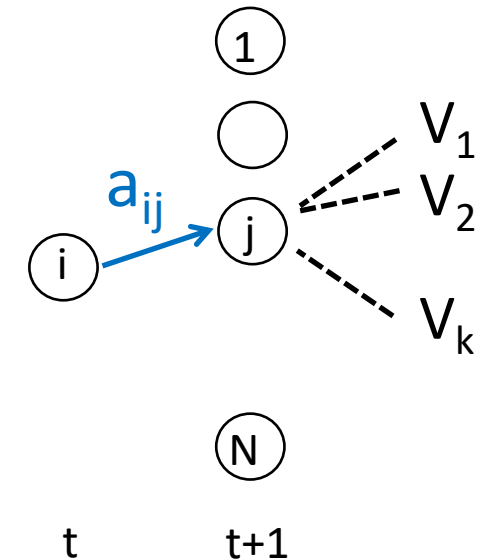
# Backward Procedure

- Define backward variable  $\beta_t(i)$ :

$$\beta_t(i) = P(O_{t+1}O_{t+2}\dots O_T | q_t = S_i, \lambda)$$

- This is the probability of the observing the **future** sequence  $O_{t+1}O_{t+2}\dots O_T$  given that the current (i.e. at time  $t$ ) state is  $S_i$ .
- It can be computed using future  $\beta_{t+1}$  as follows
  - assuming we know those future values

$$\beta_t(i) = \sum_{j=1} a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)$$



# Backward Procedure

$$\beta_t(i) = P(o_{t+1} o_{t+2} \dots o_T | \text{state at } t=i, \lambda)$$

## 1. Initialization

$$\beta_T(i) = 1, \text{ for all } i, 1 \leq i \leq N$$

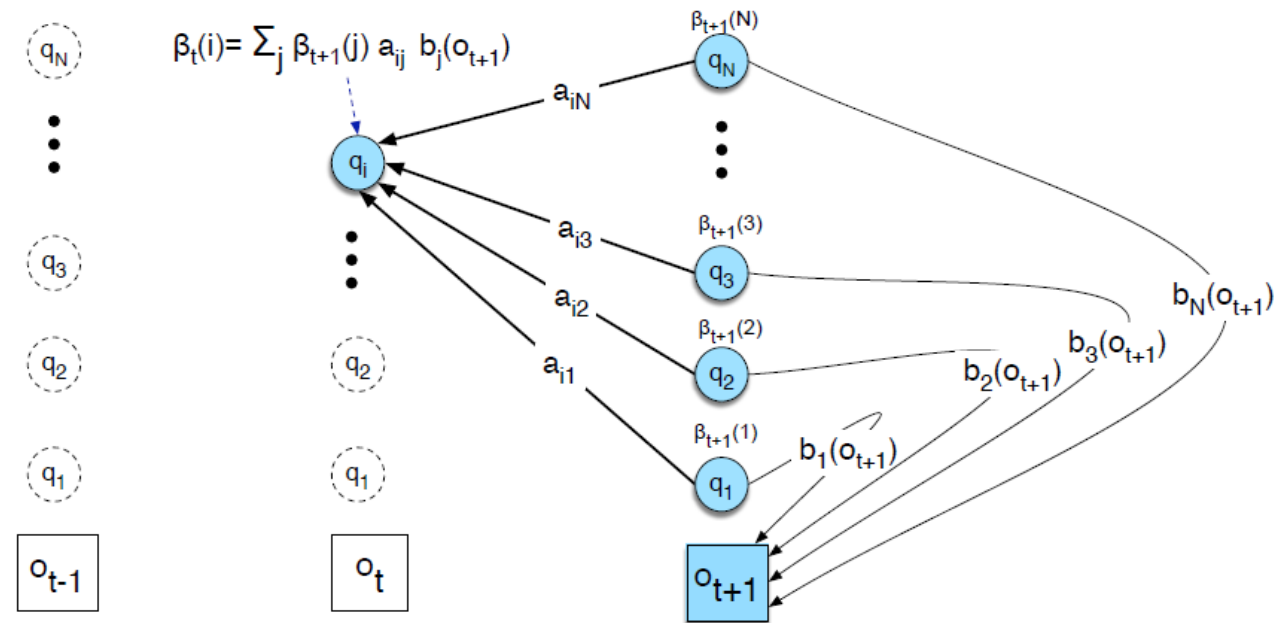
## 2. Induction

For  $t = T-1, \dots, 1$ , and for all  $i, 1 \leq i \leq N$

$$\beta_t(i) = \sum_{j=1}^N a_{ij} \beta_{t+1}(j) b_j(o_{t+1})$$

## 3. Termination

$$\text{Prob}(O | \lambda) = \sum_{i=1}^N \beta_1(i) \pi_i b_i(o_1)$$



Results of Backward & Forward Procedure  
must be the same!



# Backward Procedure

$$\beta_t(i) = P(o_{t+1} o_{t+2} \dots o_T | \text{state at } t=i, \lambda)$$

## 1. Initialization

$$\beta_T(i) = 1, \text{ for all } i, 1 \leq i \leq N$$

## 2. Induction

For  $t = T-1, \dots, 1$ , and for all  $i, 1 \leq i \leq N$

$$\beta_t(i) = \sum_{j=1}^N a_{ij} \beta_{t+1}(j) b_j(o_{t+1})$$

## 3. Termination

$$\text{Prob}(O | \lambda) = \sum_{i=1}^N \beta_1(i) \pi_i b_i(o_1)$$

$\alpha$	1	2	3						T
1									
2									
3									
4									
N									

← Fill table backward

Results of Backward & Forward Procedure  
must be the same!

# Working with Hidden Markov Models

There are three fundamental problems in HMM:

- **Evaluation Problem:**

How likely is it that HMM  $\lambda$  computed  $O$ ?

Forward or backward procedure

- **Recognition Problem:**

Does HMM  $\lambda$  recognize  $O$ ?

Viterbi Algorithm

- **Learning (= Training) Problem:**

Adjust  $\lambda$  so that  $\text{Prob}(O|\lambda)$  is locally maximized.

# Problem 2: Recognition Problem

Problem Statement: Given an observation sequence  $O=O_1O_2O_3...O_T$  and the model  $\lambda$ , what is the **optimal** state sequence  $Q = q_1q_2q_3...q_T$ ?

- “Optimal?”
  - Choose each state  $q_t$  that is **individually** most likely
  - Choose the single best **path**  $q_1q_2q_3...q_T$

i.e. to maximize  $P(Q|O,\lambda)$ ,  
equivalent to maximizing  $P(Q,O|\lambda)$

Sometimes not possible  
( $a_{ij} = 0$  for some  $i$  and  $j$ )

New Problem Statement: Given an observation sequence  $O$  and the model  $\lambda$ , what is the state sequence  $Q$  that maximizes the probability  $P(Q,O|\lambda)$ ?

## Viterbi Algorithm

- Define  $\delta_t(i)$ :

$$Q_{t-1} = q_1 q_2 \dots q_{t-1}$$

$$\delta_t(i) = \max_{Q_{t-1}} P(Q_{t-1}, q_t = S_i, O_t | \lambda)$$

- Given the model,  $\delta_t(i)$  is the best score (highest probability) along a single path, at time  $t$ , which accounts for the first  $t$  observations and ends in state  $S_i$ .
- How about  $\delta_{t+1}(j)$  for some state  $S_j$ ?  

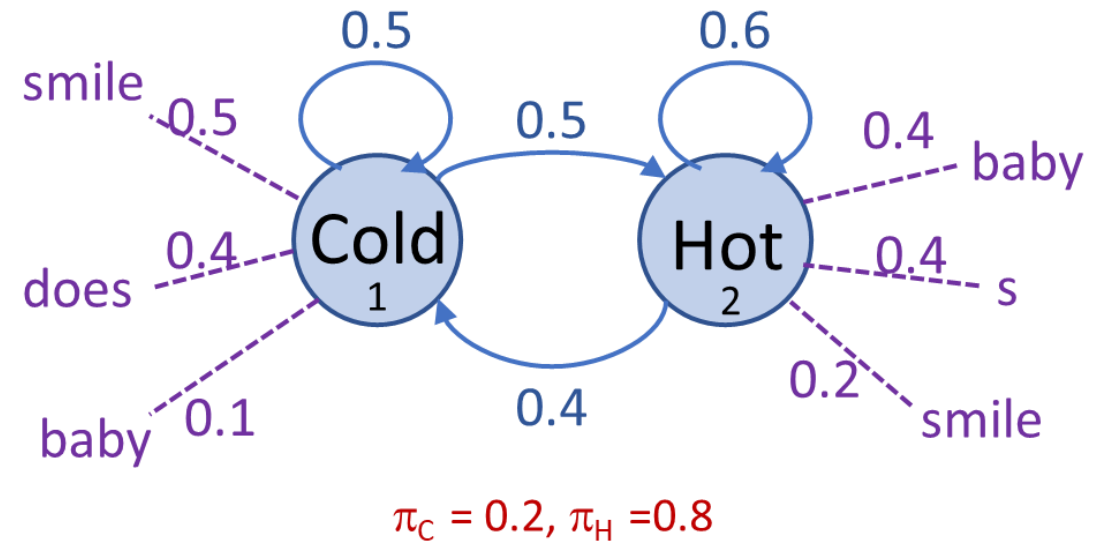
$$\delta_{t+1}(j) = [\max_i \delta_t(i) a_{ij}] b_j(O_{t+1})$$
- To actually retrieve the state sequence, we need to keep track of the argument which maximized  $\delta_{t+1}(j)$  for each  $t$  and  $j$ . We do this via the array  $\psi_t(j)$

# Viterbi Algorithm

1. Initialization:  $\delta_1(i) = \pi_i b_i(o_1)$  for all  $i$ ,  $1 \leq i \leq N$   
 $\psi_1(i) = 0$
2. Recursion: For  $t = 2, \dots, T-1$ , and for all  $j$ ,  $1 \leq j \leq N$   
 $\delta_t(j) = \max_{i=1..N} [\delta_{t-1}(i) a_{ij}] b_j(o_{t+1})$   
 $\psi_t(j) = \operatorname{argmax}_{i=1..N} [\delta_{t-1}(i) a_{ij}]$
3. Termination:  $P(Q^*, O | \lambda) = \max_{i=1..N} \delta_T(i)$   
 $q_T^* = \operatorname{argmax}_{i=1..N} \delta_T(i)$
4. Path backtracking:  
 $q_t^* = \psi_{t+1} q_{t+1}^*$  for  $t = T-1, T-2, \dots, 1$

## Example: Viterbi Algorithm

1. Initialization:  $\delta_1(i) = \pi_i b_i(o_1)$  for all  $i$ ,  $1 \leq i \leq N$   
 $\psi_1(i) = 0$
2. Recursion: For  $t = 2, \dots, T-1$ , and for all  $j$ ,  $1 \leq j \leq N$   
 $\delta_t(j) = \max_{i=1..N} [\delta_{t-1}(i) a_{ij}] b_j(o_{t+1})$   
 $\psi_t(j) = \operatorname{argmax}_{i=1..N} [\delta_{t-1}(i) a_{ij}]$
3. Termination:  $P(Q^*, O | \lambda) = \max_{i=1..N} \delta_T(i)$   
 $q_T^* = \operatorname{argmax}_{i=1..N} \delta_T(i)$
4. Path backtracking:  
 $q_t^* = \psi_{t+1} q_{t+1}^*$  for  $t = T-1, T-2, \dots, 1$

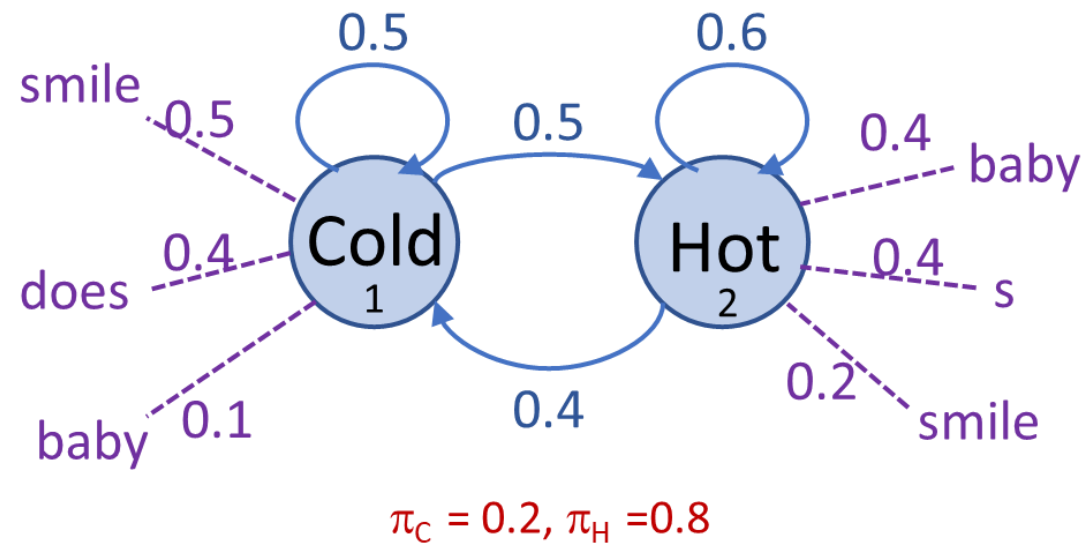


Observation  $O = \text{"baby smile s"}$  – What is the best path?

$\delta$	1	2	3
Cold			
Hot			

## Example: Viterbi Algorithm

1. Initialization:  $\delta_1(i) = \pi_i b_i(o_1)$  for all  $i$ ,  $1 \leq i \leq N$   
 $\psi_1(i) = 0$
2. Recursion: For  $t = 2, \dots, T-1$ , and for all  $j$ ,  $1 \leq j \leq N$   
 $\delta_t(j) = \max_{i=1..N} [\delta_{t-1}(i) a_{ij}] b_j(o_{t+1})$   
 $\psi_t(j) = \operatorname{argmax}_{i=1..N} [\delta_{t-1}(i) a_{ij}]$
3. Termination:  $P(Q^*, O | \lambda) = \max_{i=1..N} \delta_T(i)$   
 $q_T^* = \operatorname{argmax}_{i=1..N} \delta_T(i)$
4. Path backtracking:  
 $q_t^* = \psi_{t+1} q_{t+1}^*$  for  $t = T-1, T-2, \dots, 1$

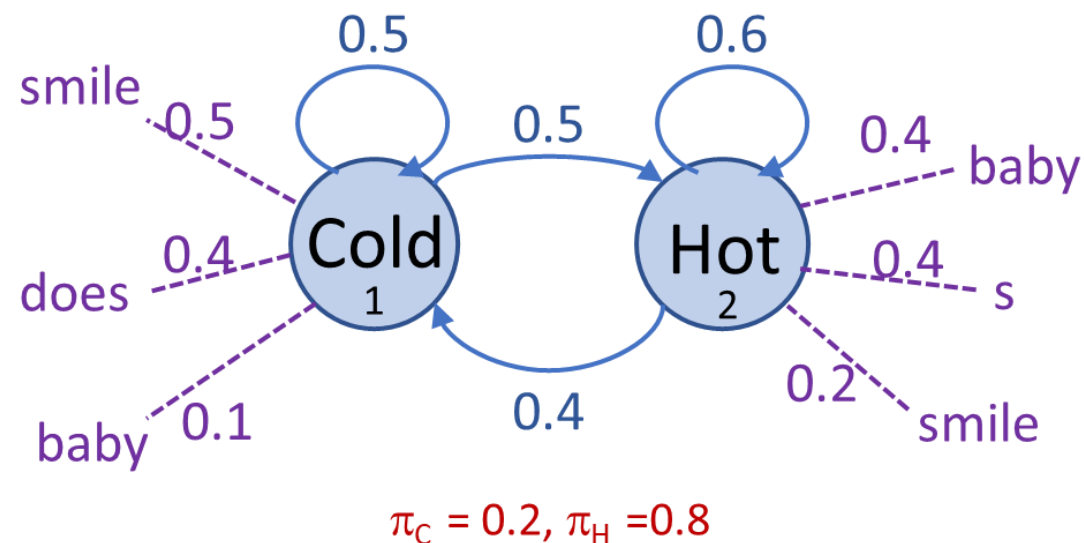


Observation O = “baby smile s” – What is the best path?

$\delta$	1	2	3
Cold	$\delta_1(C) = \pi_C b_C(\text{baby})$ $= 0.2 * 0.1 = 0.02$	$\delta_2(C)$ $= \max(\delta_1(C) \cdot a_{CC}, \delta_1(H) \cdot a_{HC}) \cdot b_C(\text{smile})$ $= \max(0.02 * 0.5, 0.32 * 0.4) * 0.5 = 0.064$	$\delta_3(C)$ $= \max(\delta_2(C) \cdot a_{CC}, \delta_2(H) \cdot a_{HC}) \cdot b_C(s)$
Hot	$\delta_1(H) = \pi_H b_H(\text{baby})$ $= 0.8 * 0.4 = 0.32$	$\delta_2(H)$ $= \max(\delta_1(C) \cdot a_{CH}, \delta_1(H) \cdot a_{HH}) \cdot b_H(\text{smile})$ $= \max(0.02 * 0.5, 0.32 * 0.6) * 0.2 = 0.038$	$\delta_3(H)$ $= \max(\delta_2(C) \cdot a_{CH}, \delta_2(H) \cdot a_{HH}) \cdot b_H(s)$

## Example: Viterbi Algorithm

1. Initialization:  $\delta_1(i) = \pi_i b_i(o_1)$  for all  $i$ ,  $1 \leq i \leq N$   
 $\psi_1(i) = 0$
2. Recursion: For  $t = 2, \dots, T-1$ , and for all  $j$ ,  $1 \leq j \leq N$   
 $\delta_t(j) = \max_{i=1..N} [\delta_{t-1}(i) a_{ij}] b_j(o_{t+1})$   
 $\psi_t(j) = \operatorname{argmax}_{i=1..N} [\delta_{t-1}(i) a_{ij}]$
3. Termination:  $P(Q^*, O | \lambda) = \max_{i=1..N} \delta_T(i)$   
 $q_T^* = \operatorname{argmax}_{i=1..N} \delta_T(i)$
4. Path backtracking:  
 $q_t^* = \psi_{t+1} q_{t+1}^*$  for  $t = T-1, T-2, \dots, 1$



Observation  $O = \text{"baby smile s"}$  – What is the best path?

$\delta$	1	2	3
Cold			
Hot			

[Code snippet](#)



# Working with Hidden Markov Models

There are three fundamental problems in HMM:

- **Evaluation Problem:**

How likely is it that HMM  $\lambda$  computed  $O$ ?

Forward or backward procedure

- **Recognition Problem:**

Does HMM  $\lambda$  recognize  $O$ ?

Viterbi Algorithm

- **Learning (= Training) Problem:**

Adjust  $\lambda$  so that  $\text{Prob}(O|\lambda)$  is locally maximized.

# Problem 3: Training HMM

Problem Statement: Given an observation sequence  $O=O_1O_2O_3...O_T$  and the model  $\lambda$ , how do we adjust  $\lambda = (A, B, \pi)$  to maximize  $P(O|\lambda)$ ?

- If we can solve this problem, then we can train a model starting from some random parameters.
- But there is no optimal way to estimate the parameters.
- One can at best use some iterative procedure to **locally** maximize the probabilities
  - Baum-Welch

# Problem 3: Training HMM

- Define

$$\xi_t(i, j) = P(q_t = S_i, q_{t+1} = S_j | O, \lambda)$$

- This is the probability of being in state  $S_i$  at time  $t$  and state  $S_j$  at time  $t + 1$ , given the observations and the model.

- Define

$$\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j) = P(q_t = S_i | O, \lambda)$$

- This is the probability of being in state  $S_i$  at time  $t$ , given the observations and the model.

# Problem 3: Training HMM

$$\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j) = P(q_t = S_i | O, \lambda)$$

$$\xi_t(i, j) = P(q_t = S_i, q_{t+1} = S_j | O, \lambda)$$

What are the sums of these two quantities over time steps  $t$  from 1 to  $T - 1$ ?

$$\sum_{t=1}^{T-1} \gamma_t(i) = \text{expected number of transitions from } S_i$$

$$\sum_{t=1}^{T-1} \xi_t(i, j) = \text{expected number of transitions from } S_i \text{ to } S_j$$

This is because they follow [Poisson binomial distribution](#).

# Problem 3: Training HMM

$$\xi_t(i, j) = P(q_t = S_i, q_{t+1} = S_j | O, \lambda)$$

$$\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j) = P(q_t = S_i | O, \lambda)$$

we can then update  $\lambda = (A, B, \pi)$  as

$$\overline{a_{ij}} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)}$$

$$\overline{b_j(v_k)} = \frac{\sum_{t=1}^{T-1} \mathbf{1}(v_k = O_t) \gamma_t(i)}{\sum_{t=1}^{T-1} \gamma_t(i)}$$

$$\overline{\pi_i} = \gamma_1(i)$$

Wait, how do we compute this?

# Problem 3: Training HMM

$$\begin{aligned}\xi_t(i, j) &= P(q_t = S_i, q_{t+1} = S_j | O, \lambda) \\ &= \frac{P(q_t = S_i, q_{t+1} = S_j, O | \lambda)}{P(O | \lambda)} \\ &= \frac{P(q_t = S_i, q_{t+1} = S_j, O | \lambda)}{\sum_i \sum_j P(q_t = S_i, q_{t+1} = S_j, O | \lambda)}\end{aligned}$$

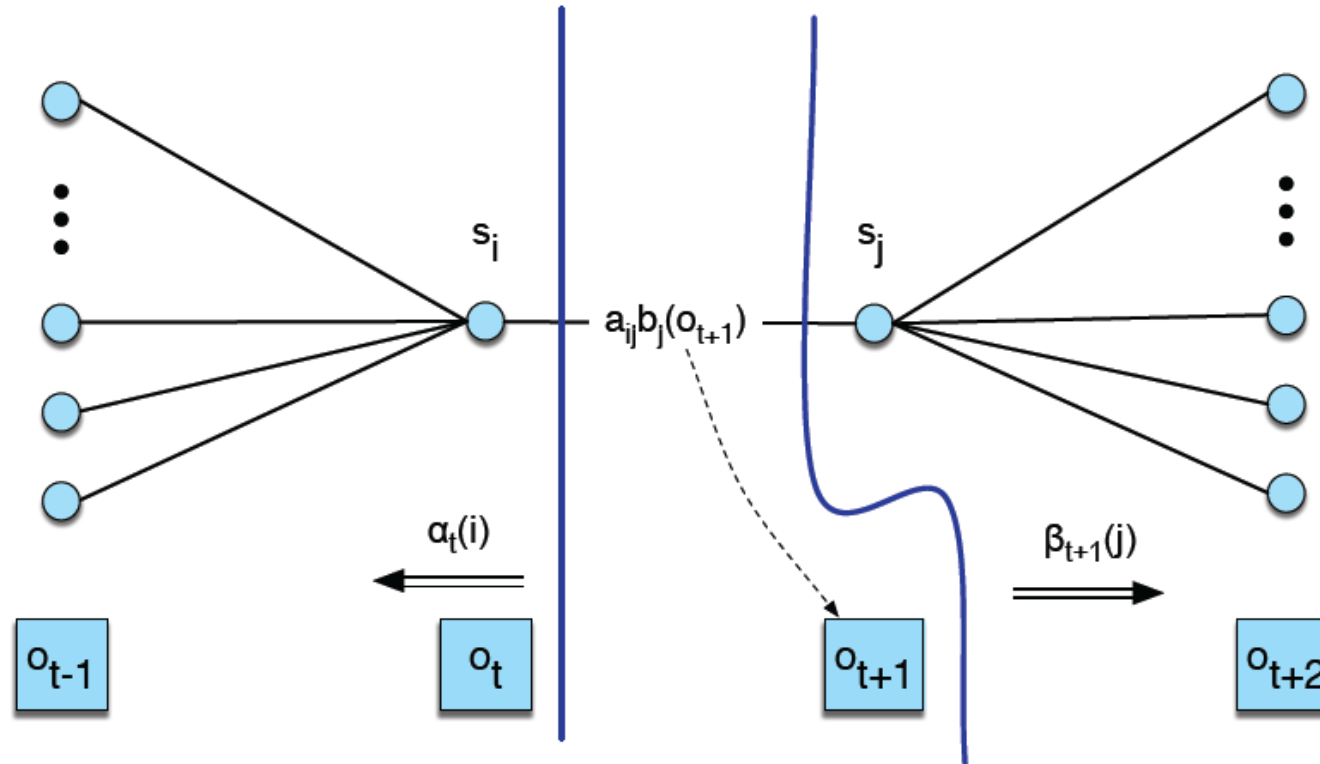
where

$$P(q_t = S_i, q_{t+1} = S_j, O | \lambda) = \alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)$$

# Recall Notations

$$\alpha_t(i) = \text{Prob}(O_1 \dots O_t \text{ and } q_t = i \mid \lambda)$$
$$\beta_t(i) = \text{Prob}(O_{t+1} \dots O_T \mid q_t = i \text{ and } \lambda)$$

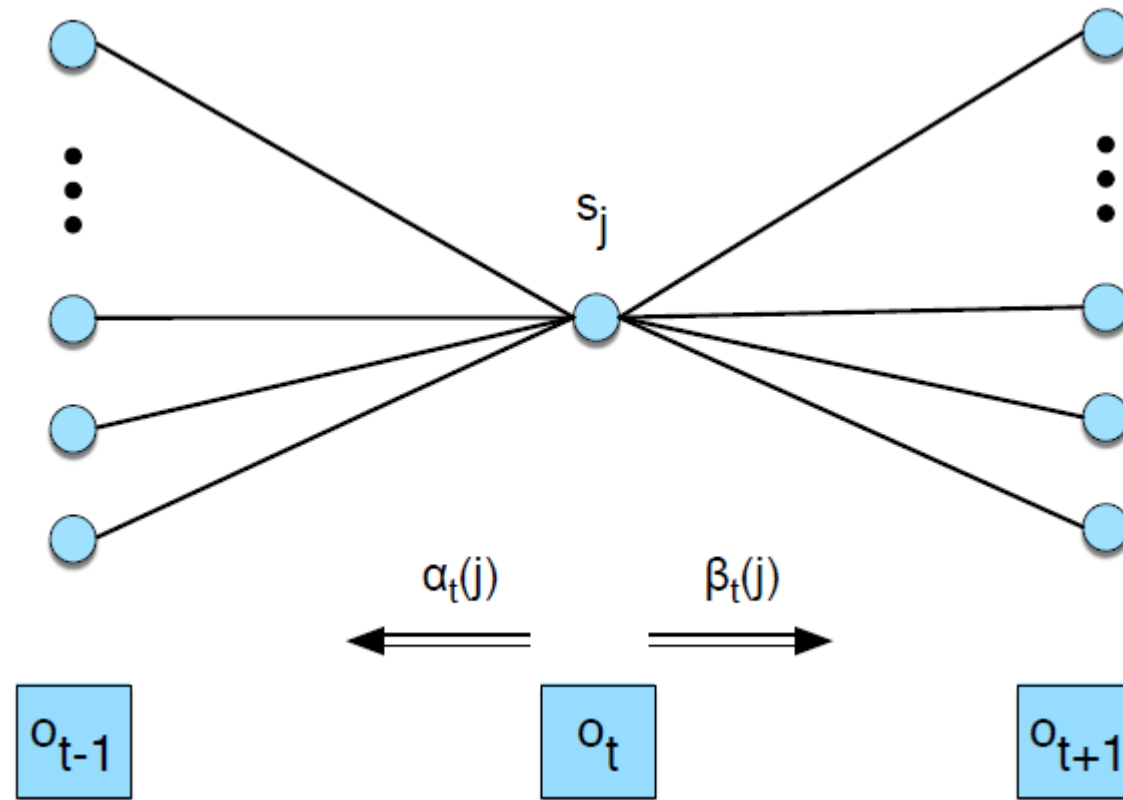
$$P(q_t = S_i, q_{t+1} = S_j, O \mid \lambda) = \alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)$$



# Recall Notations

$$\alpha_t(i) = \text{Prob}(O_1 \dots O_t \text{ and } q_t = i \mid \lambda)$$
$$\beta_t(i) = \text{Prob}(O_{t+1} \dots O_T \mid q_t = i \text{ and } \lambda)$$

$$\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j) = P(q_t = S_i \mid O, \lambda)$$





# Problem 3: Training HMM

$$\xi_t(i, j) = P(q_t = S_i, q_{t+1} = S_j | O, \lambda)$$

$$\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j) = P(q_t = S_i | O, \lambda)$$

$$\alpha_t(i) = \text{Prob}(O_1 \dots O_t \text{ and } q_t = i \mid \lambda)$$

$$\beta_t(i) = \text{Prob}(O_{t+1} \dots O_T \mid q_t = i \text{ and } \lambda)$$

$$\xi_t(i, j) = \alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j) / \text{Prob}(O \mid \lambda)$$

$$= \alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j) / \sum_{i=1}^n \sum_{j=1}^n \alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)$$

$$\gamma_t(i) = \alpha_t(i) \beta_t(i) / \text{Prob}(O \mid \lambda)$$

$$= \alpha_t(i) \beta_t(i) / \sum_{i=1}^n \alpha_t(i) \beta_t(i)$$

# Baum-Welch Re-estimation

1.  $\bar{\pi}_i = \gamma_1(i)$
2.  $\bar{a}_{ij} = \sum_{t=1}^T \xi_t(i,j) / \sum_{t=1}^T \gamma_t(i)$
3.  $\bar{b}_j(k) = \sum_{t=1, o_t=v_k}^T \gamma_t(j) / \sum_{t=1}^T \gamma_t(j)$

If  $\text{Prob}(O \mid \bar{\lambda}) > \text{Prob}(O \mid \lambda)$

Re-estimate  $\bar{\lambda}$  until (local) maximum is reached

# Problem3: Training HMMs

- No analytic solution ( = no “simply plug in” solution)
- Use iterative algorithm that learns to represent training data better and better

“Baum-Welch Reestimation Algorithm”

Input: HMM  $\lambda$ , training sequence  $O$

- 1) Initialization: Guess all probabilities to be uniform
- 2) Repeat until no better HMM  $\lambda$  can be found:

Update all probabilities according to equations on next slide

Output: HMM  $\lambda^*$

# Learning Outcomes

- Understand FSM, Deterministic FSM, Markov Network(Chains), HMM
  - The “Markov” Property
  - The “Hidden”
- Know the three fundamental problems in HMM, and how to solve them
  - Forward procedure
  - Viterbi Algorithm
  - Baum-Welch Reestimation Algorithm – using the forward-backward quantities