

# CS344: Introduction to Artificial Intelligence (associated lab: CS386)

Pushpak Bhattacharyya  
CSE Dept.,  
IIT Bombay

Lecture 9: Viterbi; forward and backward  
probabilities

25<sup>th</sup> Jan, 2011

# HMM Definition

- Set of states:  $S$  where  $|S|=N$
- Start state  $S_0$  /\* $P(S_0)=1$ \*/
- Output Alphabet:  $O$  where  $|O|=M$
- Transition Probabilities:  $A = \{a_{ij}\}$  /\*state  $i$  to state  $j$ \*/
- Emission Probabilities :  $B = \{b_j(o_k)\}$  /\*prob. of emitting or absorbing  $o_k$  from state  $j$ \*/
- Initial State Probabilities:  $\Pi = \{p_1, p_2, p_3, \dots, p_N\}$
- Each  $p_i = P(o_0 = \varepsilon, S_i | S_0)$

# Markov Processes

- Properties

- Limited Horizon: Given previous  $t$  states, a state  $i$ , is independent of preceding  $0$  to  $t-k+1$  states.
  - $P(X_t=i/X_{t-1}, X_{t-2}, \dots, X_0) = P(X_t=i/X_{t-1}, X_{t-2}, \dots, X_{t-k})$
  - Order  $k$  Markov process
- Time invariance: (shown for  $k=1$ )
  - $P(X_t=i/X_{t-1}=j) = P(X_1=i/X_0=j) \dots = P(X_n=i/X_{n-1}=j)$

# Three basic problems (contd.)

- Problem 1: Likelihood of a sequence
  - Forward Procedure
  - Backward Procedure
- Problem 2: Best state sequence
  - Viterbi Algorithm
- Problem 3: Re-estimation
  - Baum-Welch ( Forward-Backward Algorithm )

# Probabilistic Inference

- O: Observation Sequence
- S: State Sequence
- Given O find  $S^*$  where  $S^* = \arg \max_S p(S / O)$  called Probabilistic Inference
- Infer “Hidden” from “Observed”
- How is this inference different from logical inference based on propositional or predicate calculus?

# Essentials of Hidden Markov Model

1. Markov + Naive Bayes
2. Uses both transition and observation probability

$$p(S_k \rightarrow^{O_k} S_{k+1}) = p(O_k / S_k) p(S_{k+1} / S_k)$$

3. Effectively makes Hidden Markov Model a Finite State Machine (FSM) with probability

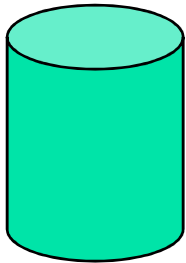
# Probability of Observation Sequence

$$\begin{aligned} p(O) &= \sum_S p(O, S) \\ &= \sum_S p(S) p(O / S) \end{aligned}$$

- Without any restriction,
  - Search space size =  $|S|^{|O|}$

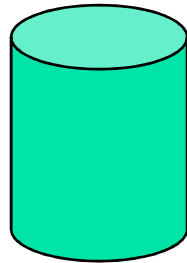
# Continuing with the Urn example

Colored Ball choosing



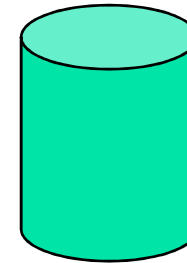
Urn 1

# of Red = 30  
# of Green = 50  
# of Blue = 20



Urn 2

# of Red = 10  
# of Green = 40  
# of Blue = 50



Urn 3

# of Red = 60  
# of Green = 10  
# of Blue = 30



# Example (contd.)

Transition Probability

	$U_1$	$U_2$	$U_3$
$U_1$	0.1	0.4	0.5
$U_2$	0.6	0.2	0.2
$U_3$	0.3	0.4	0.3

Given :

Observation/output Probability

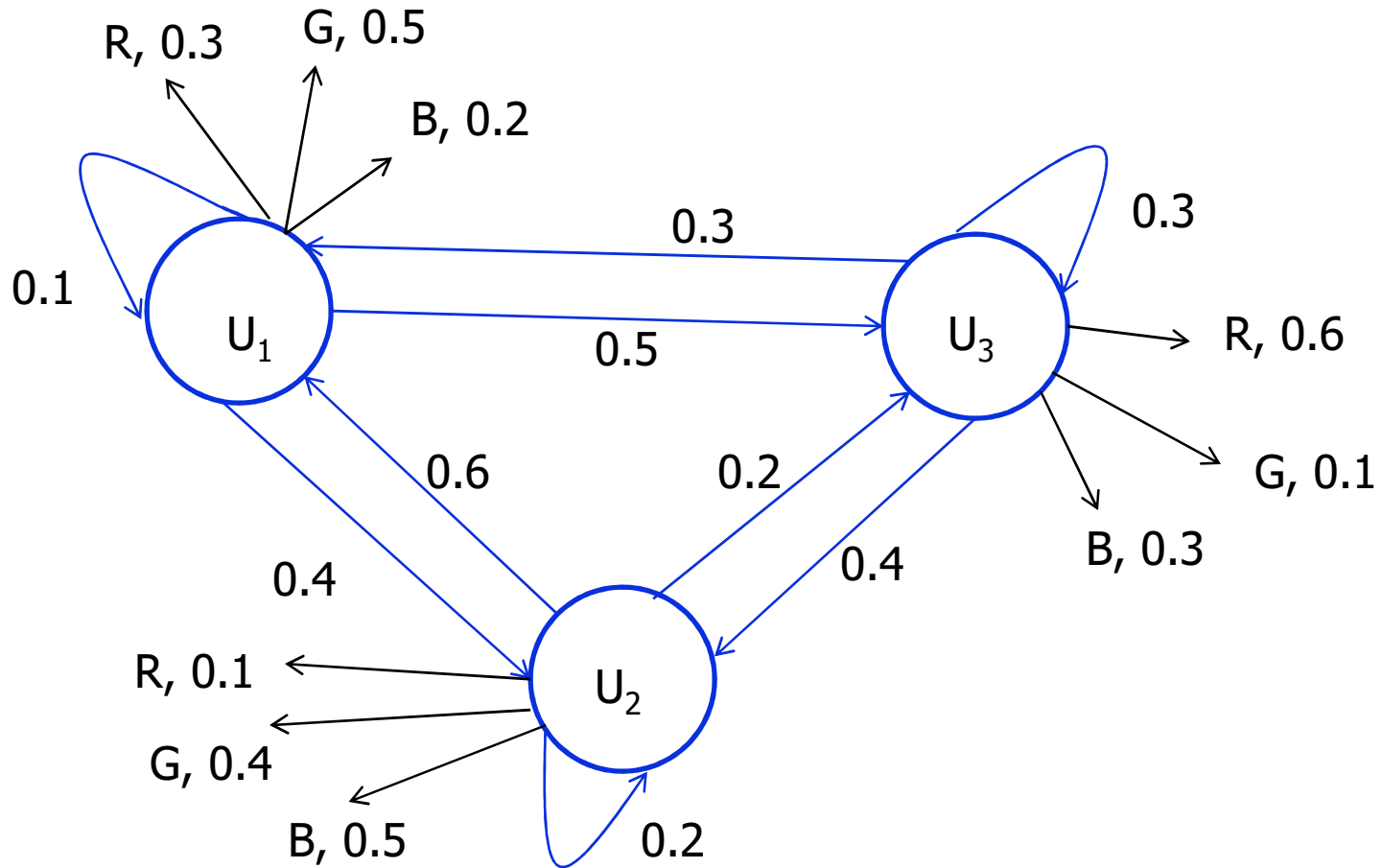
	R	G	B
$U_1$	0.3	0.5	0.2
$U_2$	0.1	0.4	0.5
$U_3$	0.6	0.1	0.3

and

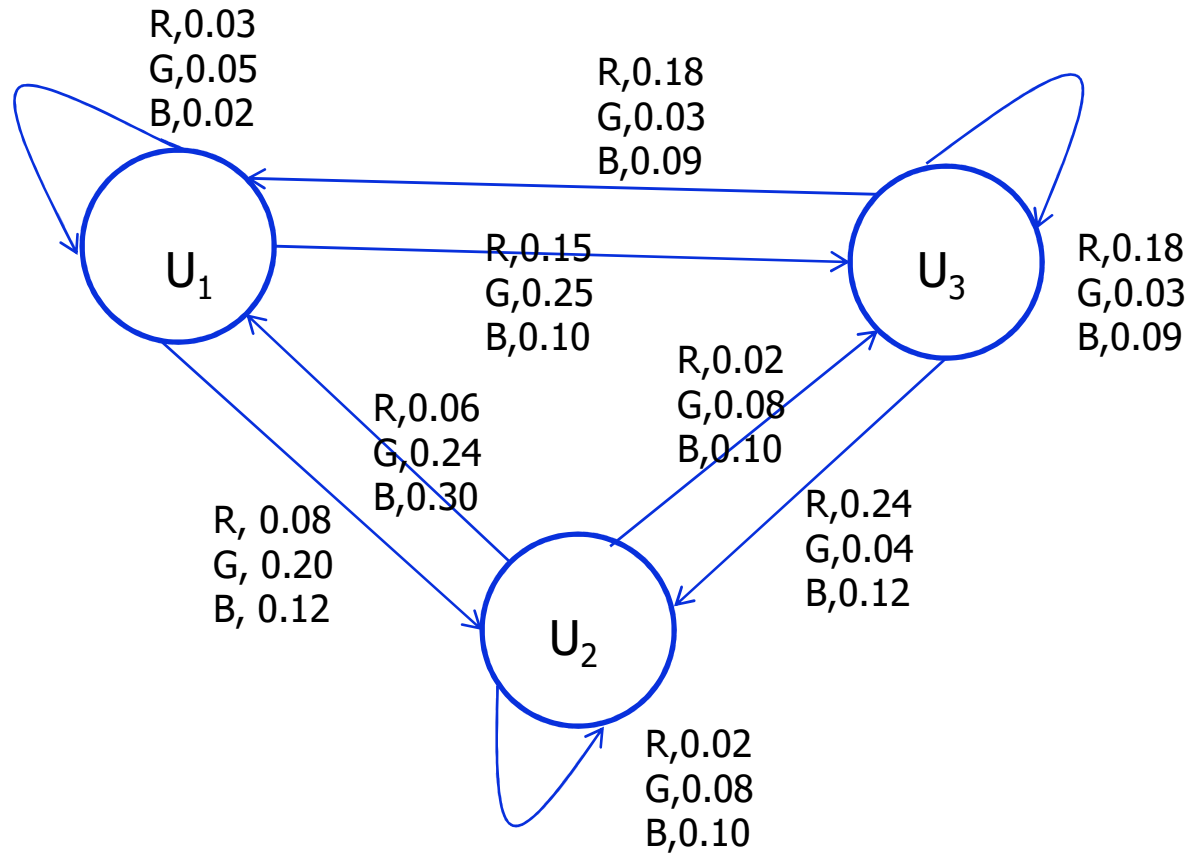
Observation : RRGGBRGR

What is the corresponding state sequence ?

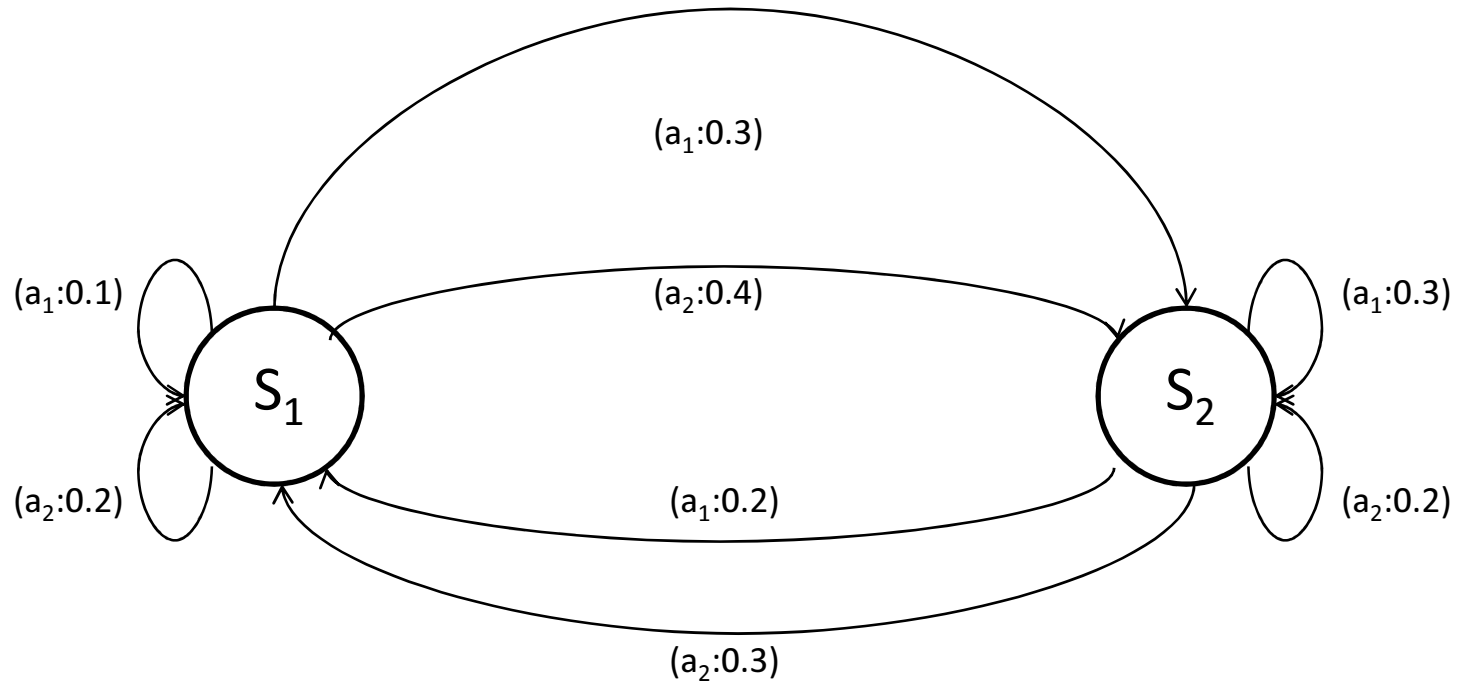
# Diagrammatic representation (1/2)



# Diagrammatic representation (2/2)



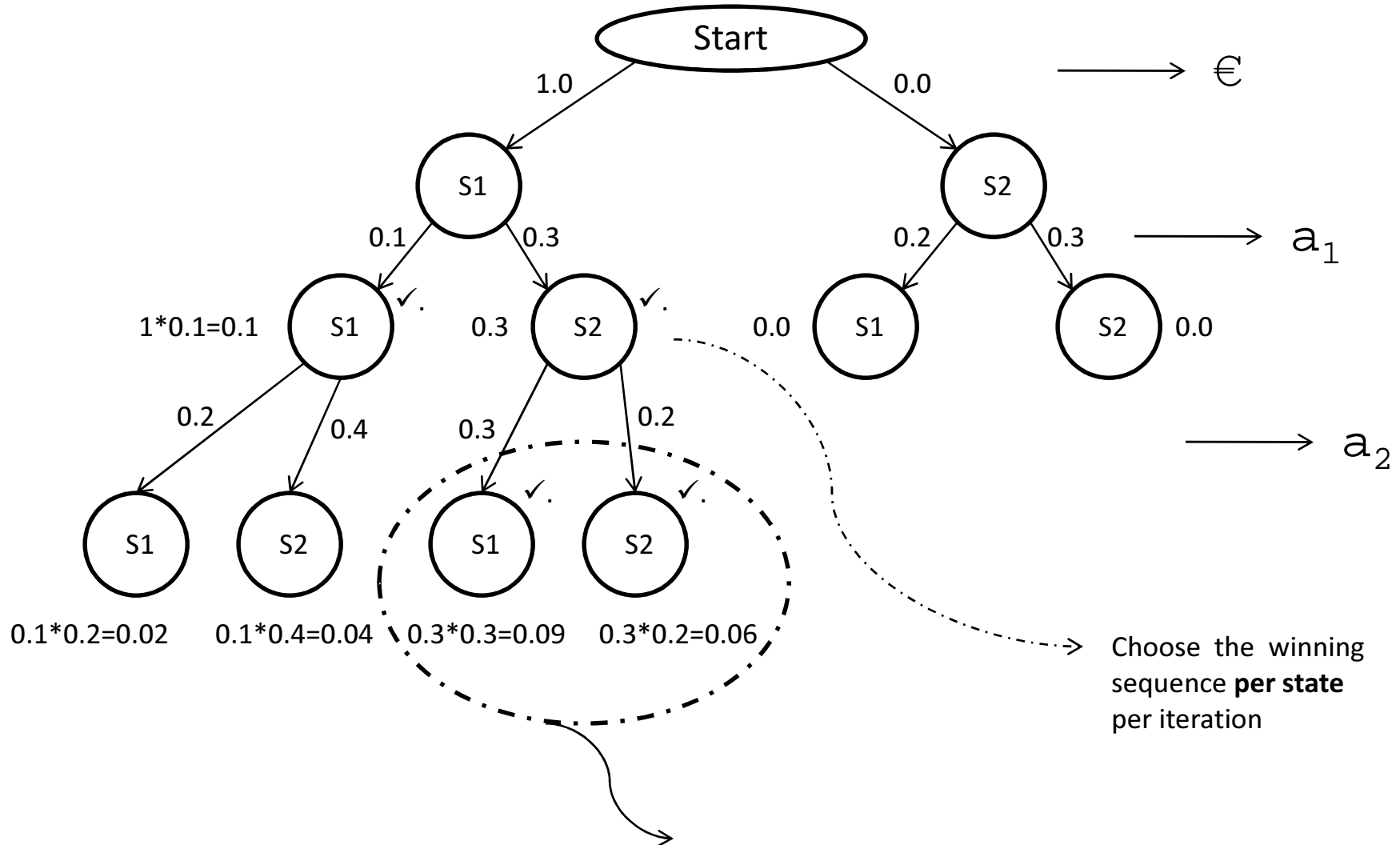
# Probabilistic FSM



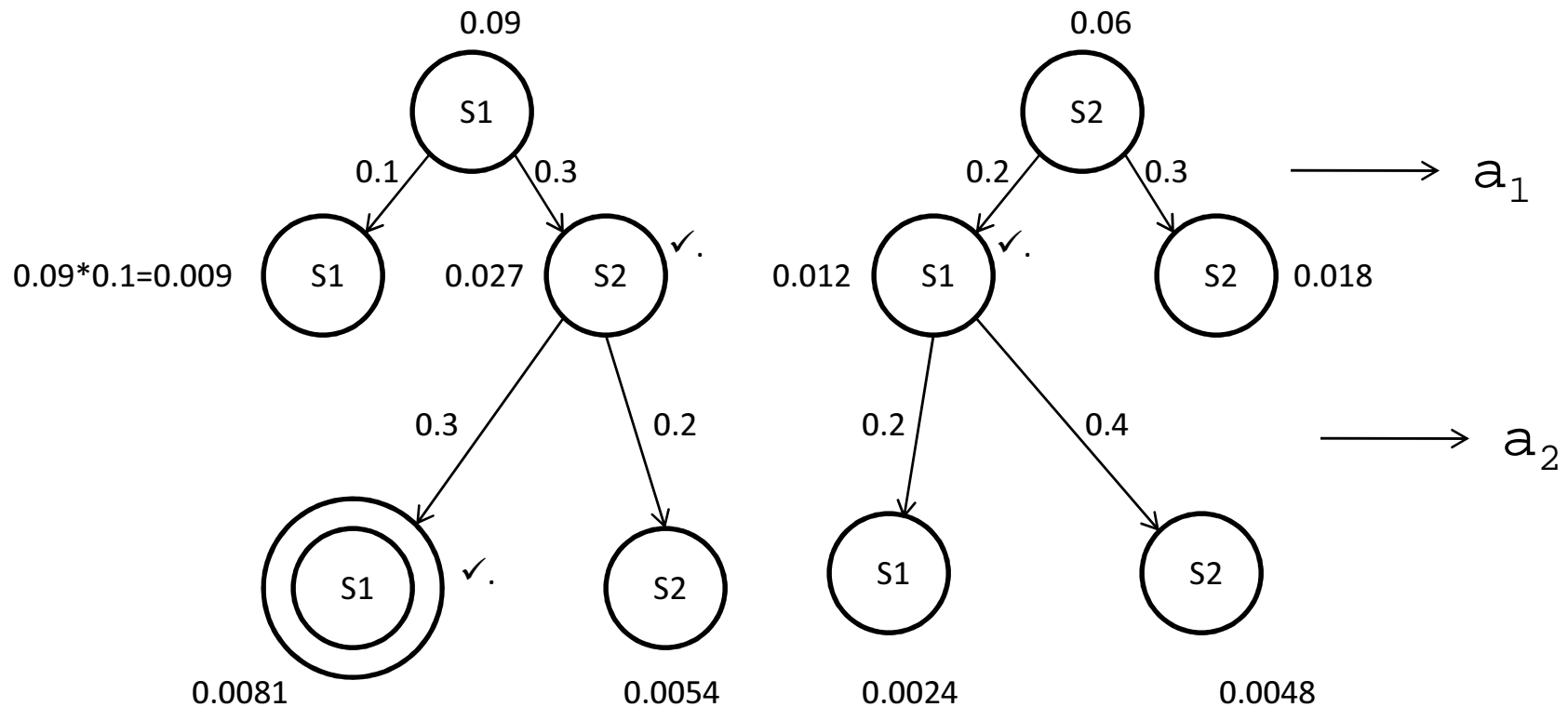
The question here is:

“what is the most likely state sequence given the output sequence seen”

# Developing the tree



# Tree structure contd...



The problem being addressed by this tree is  $S^* = \arg \max_s P(S | a_1 - a_2 - a_1 - a_2, \mu)$

$a_1 - a_2 - a_1 - a_2$  is the output sequence and  $\mu$  the model or the machine

# Tabular representation of the tree

Ending state \ Latest symbol observed	€	$a_1$	$a_2$	$a_1$	$a_2$
$S_1$	1.0	$(1.0*0.1, 0.0*0.2) = (\mathbf{0.1}, 0.0)$	$(0.02, \mathbf{0.09})$	$(0.009, \mathbf{0.012})$	$(0.0024, \mathbf{0.0081})$
$S_2$	0.0	$(1.0*0.3, 0.0*0.3) = (\mathbf{0.3}, 0.0)$	$(0.04, \mathbf{0.06})$	$(\mathbf{0.027}, 0.018)$	$(0.0048, 0.0054)$

**Note:** Every cell records the winning probability ending in that state

Final winner

The bold faced values in each cell shows the sequence probability ending in that state. Going backward from final winner sequence which ends in state  $S_2$  (indicated by the 2<sup>nd</sup> tuple), we recover the sequence.

# Algorithm

*(following James Alan, Natural Language Understanding (2<sup>nd</sup> edition), Benjamin Cummins (pub.), 1995)*

Given:

1. The HMM, which means:
  - a. Start State:  $S_1$
  - b. Alphabet:  $A = \{a_1, a_2, \dots, a_p\}$
  - c. Set of States:  $S = \{S_1, S_2, \dots, S_n\}$
  - d. Transition probability  $P(S_i \xrightarrow{a^k} S_j) \quad \forall i, j, k$   
which is equal to  $P(S_j, a^k | S_i)$
2. The output string  $a_1 a_2 \dots a_T$

To find:

The most likely sequence of states  $C_1 C_2 \dots C_T$  which produces the given output sequence, *i.e.*,  $C_1 C_2 \dots C_T = \arg \max_C [P(C | a_1, a_2, \dots, a_T, \mu)]$



# Algorithm contd...

## Data Structure:

1. A  $N \times T$  array called SEQSCORE to maintain the winner sequence always ( $N = \# \text{states}$ ,  $T = \text{length of o/p sequence}$ )
2. Another  $N \times T$  array called BACKPTR to recover the path.

## Three distinct steps in the Viterbi implementation

1. Initialization
2. Iteration
3. Sequence Identification

# 1. Initialization

SEQSCORE(1,1)=1.0

BACKPTR(1,1)=0

For(i=2 to N) do

    SEQSCORE(i,1)=0.0

[expressing the fact that first state  
is  $S_1$ ]

# 2. Iteration

For(t=2 to T) do

    For(i=1 to N) do

        SEQSCORE(i,t) =  $\text{Max}_{(j=1,N)}$

        [  $\text{SEQSCORE}(j, (t - 1)) * P(S_j \xrightarrow{a_k} S_i)$  ]

        BACKPTR(i,t) = index  $j$  that gives the MAX above

### 3. Seq. Identification

$C(T) = i$  that maximizes  $SEQSCORE(i,T)$

For  $i$  from  $(T-1)$  to  $1$  do

$C(i) = BACKPTR[C(i+1),(i+1)]$

#### Optimizations possible:

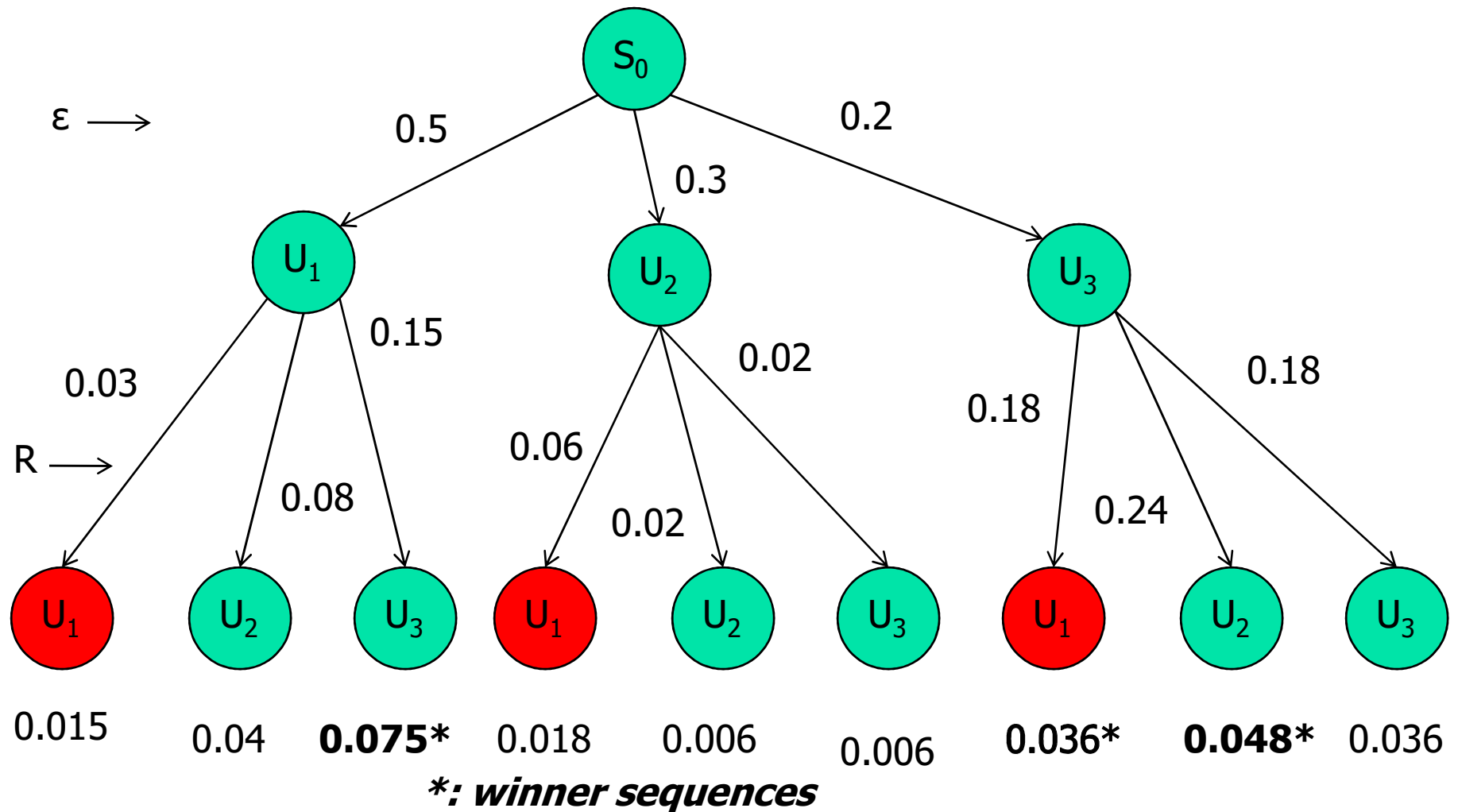
1.  $BACKPTR$  can be  $1*T$
2.  $SEQSCORE$  can be  $T*2$

**Homework:-** Compare this with  $A^*$ , Beam Search [Homework]

Reason for this comparison:

Both of them work for finding and recovering sequence

# Viterbi Algorithm for the Urn problem (first two symbols)



# Markov process of order > 1 (say 2)

	$O_0$	$O_1$	$O_2$	$O_3$	$O_4$	$O_5$	$O_6$	$O_7$	$O_8$	
Obs:	$\epsilon$	R	R	G	G	B	R	G	R	
State:	$S_0$	$S_1$	$S_2$	$S_3$	$S_4$	$S_5$	$S_6$	$S_7$	$S_8$	$S_9$

Same theory works

$P(S).P(O|S)$

$$\begin{aligned}
 &= P(O_0|S_0).P(S_1|S_0). \\
 &\quad [P(O_1|S_1).P(S_2|S_1S_0)]. \\
 &\quad [P(O_2|S_2).P(S_3|S_2S_1)]. \\
 &\quad [P(O_3|S_3).P(S_4|S_3S_2)]. \\
 &\quad [P(O_4|S_4).P(S_5|S_4S_3)]. \\
 &\quad [P(O_5|S_5).P(S_6|S_5S_4)]. \\
 &\quad [P(O_6|S_6).P(S_7|S_6S_5)]. \\
 &\quad [P(O_7|S_7).P(S_8|S_7S_6)]. \\
 &\quad [P(O_8|S_8).P(S_9|S_8S_7)].
 \end{aligned}$$

We introduce the states  $S_0$  and  $S_9$  as initial and final states respectively.

After  $S_8$  the next state is  $S_9$  with probability 1, i.e.,  $P(S_9|S_8S_7)=1$

$O_0$  is  $\epsilon$ -transition

# Adjustments

- Transition probability table will have tuples on rows and states on columns
- Output probability table will remain the same
- In the Viterbi tree, the Markov process will take effect from the 3<sup>rd</sup> input symbol ( $\epsilon RR$ )
- There will be 27 leaves, out of which **only 9 will remain**
- Sequences ending in **same tuples** will be compared
  - Instead of  $U_1, U_2$  and  $U_3$
  - $U_1U_1, U_1U_2, U_1U_3, U_2U_1, U_2U_2, U_2U_3, U_3U_1, U_3U_2, U_3U_3$

# Forward and Backward Probability Calculation

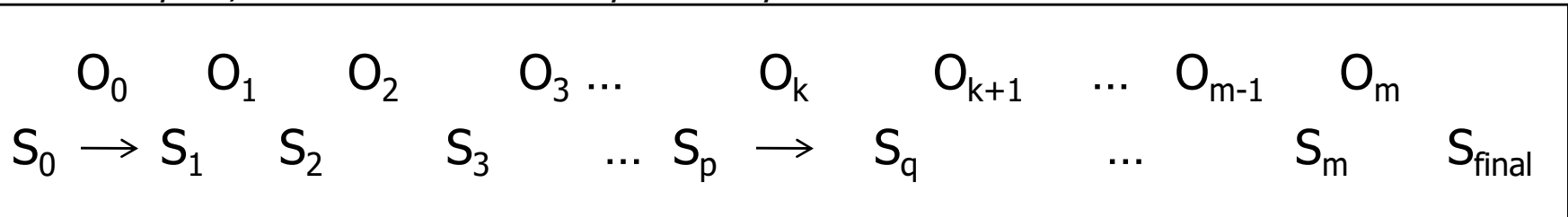
# Forward probability $F(k,i)$

- Define  $F(k,i)$  = Probability of being in state  $S_i$  having seen  $o_0o_1o_2\dots o_k$
- $F(k,i) = P(o_0o_1o_2\dots o_k, S_i)$
- With  $m$  as the length of the observed sequence
- $P(\text{observed sequence}) = P(o_0o_1o_2\dots o_m)$   
 $= \sum_{p=0,N} P(o_0o_1o_2\dots o_m, S_p)$   
 $= \sum_{p=0,N} F(m, p)$



# Forward probability (contd.)

$$\begin{aligned}
 F(k, q) &= P(o_0 o_1 o_2 \dots o_k, S_q) \\
 &= P(o_0 o_1 o_2 \dots o_k, S_q) \\
 &= P(o_0 o_1 o_2 \dots o_{k-1}, o_k, S_q) \\
 &= \sum_{p=0, N} P(o_0 o_1 o_2 \dots o_{k-1}, S_p, o_k, S_q) \\
 &= \sum_{p=0, N} P(o_0 o_1 o_2 \dots o_{k-1}, S_p) \cdot \\
 &\quad P(o_k, S_q / o_0 o_1 o_2 \dots o_{k-1}, S_p) \\
 &= \sum_{p=0, N} F(k-1, p) \cdot P(o_k, S_q / S_p) \\
 &= \sum_{p=0, N} F(k-1, p) \cdot P(S_p \xrightarrow{o_k} S_q)
 \end{aligned}$$



# Backward probability $B(k,i)$

- Define  $B(k,i)$  = Probability of seeing  $O_k O_{k+1} O_{k+2} \dots O_m$  given that the state was  $S_i$
- $B(k,i) = P(O_k O_{k+1} O_{k+2} \dots O_m \mid S_i)$
- With  $m$  as the length of the observed sequence
- $P(\text{observed sequence}) = P(O_0 O_1 O_2 \dots O_m)$   
 $= P(O_0 O_1 O_2 \dots O_m \mid S_0)$   
 $= B(0,0)$

# Backward probability (contd.)

$$\begin{aligned}
 & B(k, p) \\
 &= P(o_k o_{k+1} o_{k+2} \dots o_m \mid S_p) \\
 &= P(o_{k+1} o_{k+2} \dots o_m, o_k \mid S_p) \\
 &= \sum_{q=0, N} P(o_{k+1} o_{k+2} \dots o_m, o_k, S_q \mid S_p) \\
 &= \sum_{q=0, N} P(o_k, S_q \mid S_p) \\
 &\quad P(o_{k+1} o_{k+2} \dots o_m \mid o_k, S_q, S_p) \\
 &= \sum_{q=0, N} P(o_{k+1} o_{k+2} \dots o_m \mid S_q) \cdot P(o_k, S_q \mid S_p) \\
 &= \sum_{q=0, N} B(k+1, q) \cdot P(S_p \xrightarrow{o_k} S_q)
 \end{aligned}$$

