

Hidden Markov Models

Manuela Veloso

Co-instructor: Pradeep Ravikumar

Machine Learning 10-701

Some slides courtesy of Andrew Moore and Eric Xing

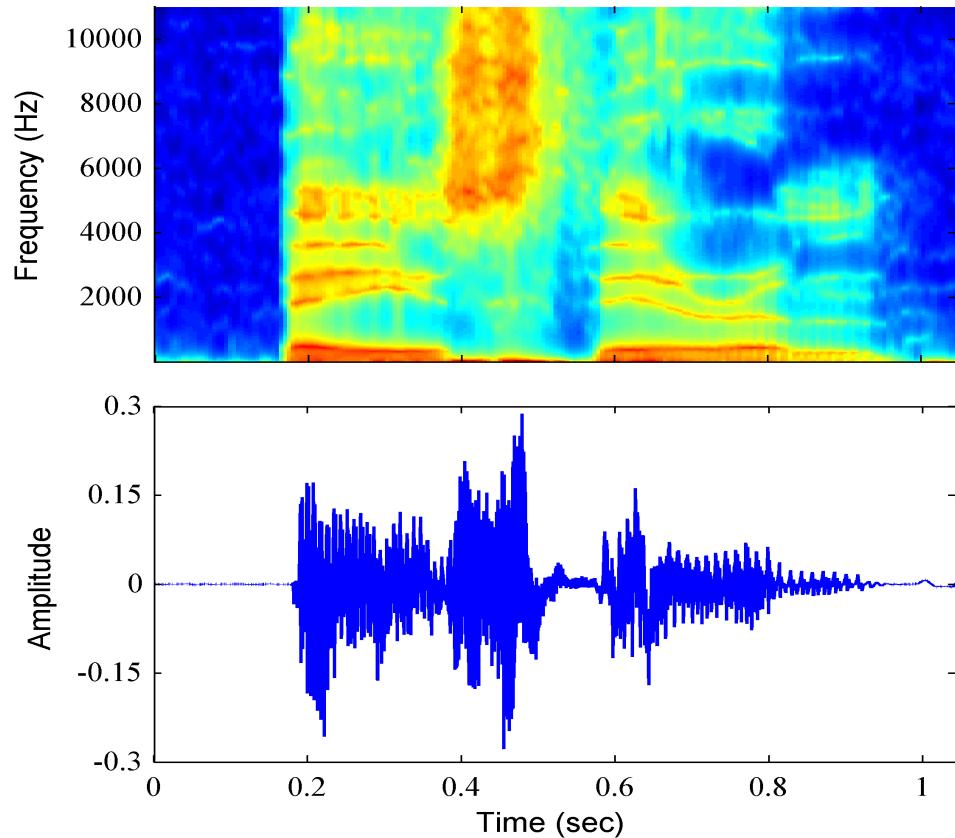


MACHINE LEARNING DEPARTMENT

Carnegie Mellon.
School of Computer Science

Sequential data

- So far we assumed independent, identically distributed data, $\{X_i\}_{i=1}^n \stackrel{iid}{\sim} p(\mathbf{X})$
- Sequential data
 - Time-series data
E.g. Speech



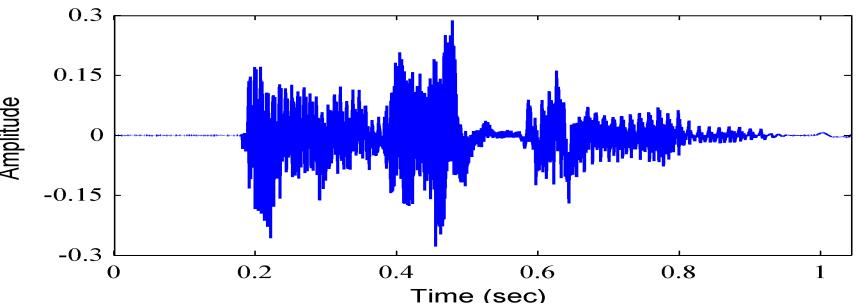
i.i.d. to sequential data

- So far we assumed independent, identically distributed data, $\{X_i\}_{i=1}^n \stackrel{iid}{\sim} p(\mathbf{X})$

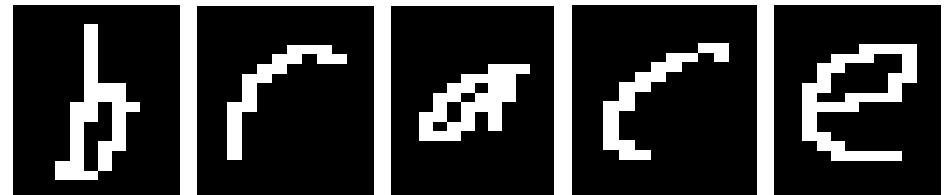
- Sequential data

- Time-series data

- E.g. Speech



- Characters in a sentence



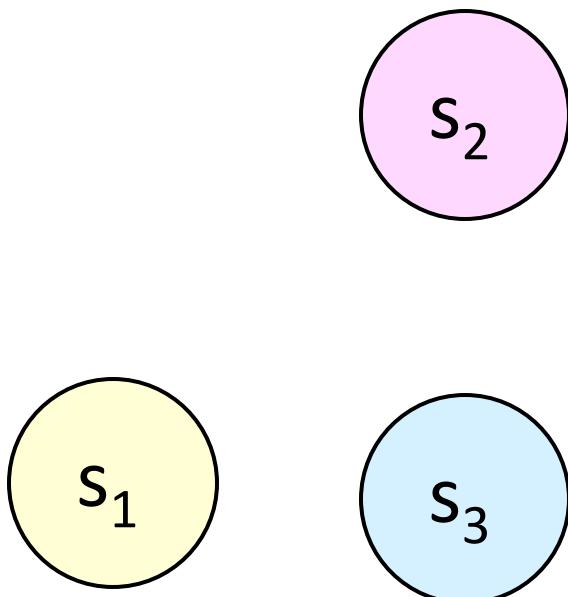
- Mobile robots



A Markov System

N states, called $s_1, s_2 \dots s_N$

Discrete timesteps, $t=0, t=1, \dots$



$N = 3$

$t=0$

A Markov System

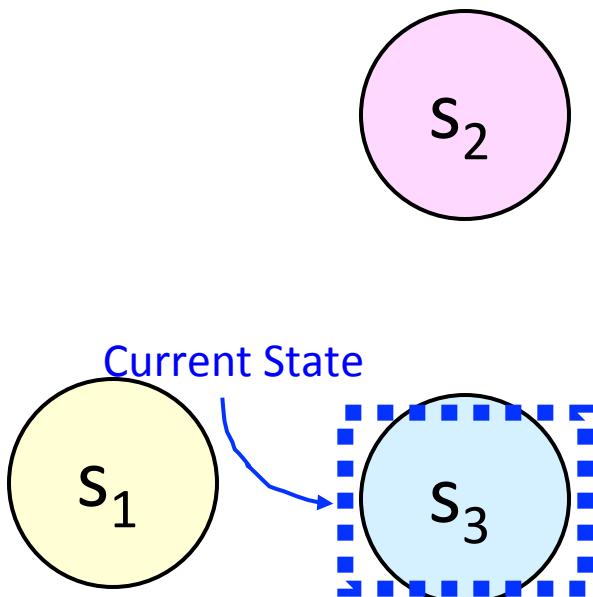
N states, called $s_1, s_2 \dots s_N$

Discrete timesteps, $t=0, t=1, \dots$

On timestep t , the system is in exactly one of the available states.

Call it q_t

Note: $q_t \in \{s_1, s_2 \dots s_N\}$

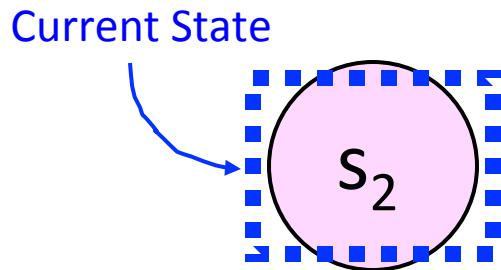


$$N = 3$$

$$t=0$$

$$q_t = q_0 = s_3$$

A Markov System



$N = 3$

$t=1$

$q_t = q_1 = s_2$

N states, called $s_1, s_2 \dots s_N$

Discrete timesteps, $t=0, t=1, \dots$

On timestep t , the system is in exactly one of the available states.

Call it q_t

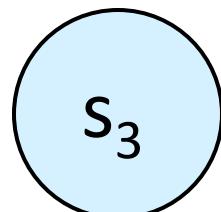
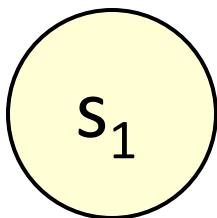
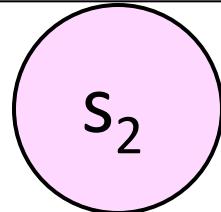
Note: $q_t \in \{s_1, s_2 \dots s_N\}$

After a timestep, the next state is chosen randomly according to a transition probability.

A Markov System

$$\begin{aligned} P(q_{t+1}=s_1 | q_t=s_2) &= 1/2 \\ P(q_{t+1}=s_2 | q_t=s_2) &= 1/2 \\ P(q_{t+1}=s_3 | q_t=s_2) &= 0 \end{aligned}$$

$$\begin{aligned} P(q_{t+1}=s_1 | q_t=s_1) &= 0 \\ P(q_{t+1}=s_2 | q_t=s_1) &= 0 \\ P(q_{t+1}=s_3 | q_t=s_1) &= 1 \end{aligned}$$



$N = 3$

$t=1$

$q_t=q_1=s_2$

$$\begin{aligned} P(q_{t+1}=s_1 | q_t=s_3) &= 1/3 \\ P(q_{t+1}=s_2 | q_t=s_3) &= 2/3 \\ P(q_{t+1}=s_3 | q_t=s_3) &= 0 \end{aligned}$$

N states, called $s_1, s_2 \dots s_N$

Discrete timesteps, $t=0, t=1, \dots$

On timestep t , the system is in exactly one of the available states.

Call it q_t

Note: $q_t \in \{s_1, s_2 \dots s_N\}$

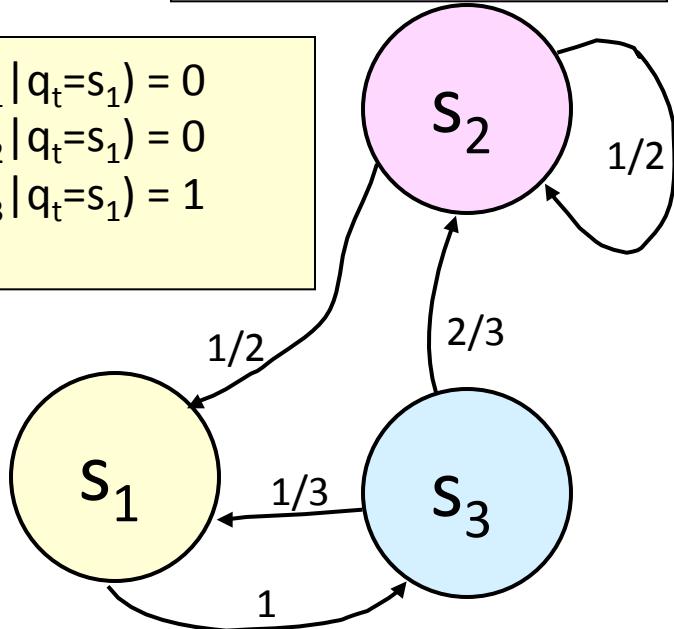
After a timestep, the next state is chosen randomly according to a transition probability.

The current state determines the probability distribution for the next state.

A Markov System

$$\begin{aligned} P(q_{t+1}=s_1 | q_t=s_2) &= 1/2 \\ P(q_{t+1}=s_2 | q_t=s_2) &= 1/2 \\ P(q_{t+1}=s_3 | q_t=s_2) &= 0 \end{aligned}$$

$$\begin{aligned} P(q_{t+1}=s_1 | q_t=s_1) &= 0 \\ P(q_{t+1}=s_2 | q_t=s_1) &= 0 \\ P(q_{t+1}=s_3 | q_t=s_1) &= 1 \end{aligned}$$



$N = 3$

$t=1$

$q_t=q_1=s_2$

$$\begin{aligned} P(q_{t+1}=s_1 | q_t=s_3) &= 1/3 \\ P(q_{t+1}=s_2 | q_t=s_3) &= 2/3 \\ P(q_{t+1}=s_3 | q_t=s_3) &= 0 \end{aligned}$$

N states, called $s_1, s_2 \dots s_N$

Discrete timesteps, $t=0, t=1, \dots$

On timestep t , the system is in exactly one of the available states.

Call it q_t

Note: $q_t \in \{s_1, s_2 \dots s_N\}$

After a timestep, the next state is chosen randomly according to a transition probability.

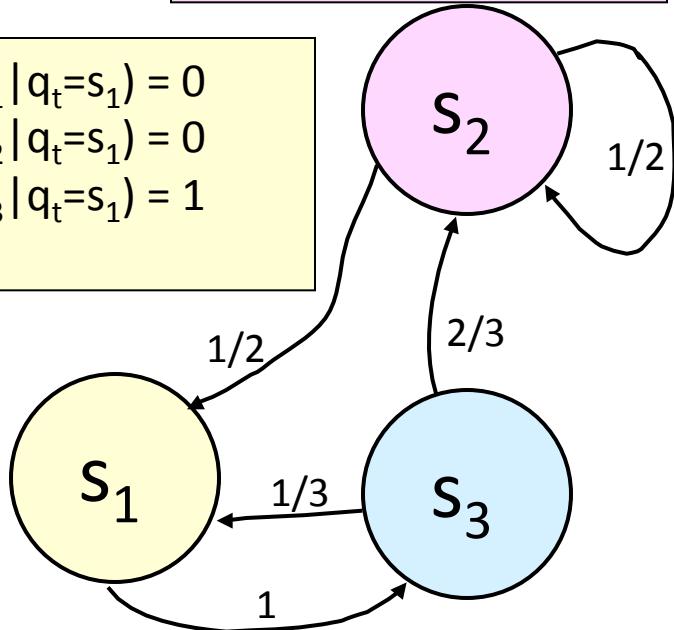
The current state determines the probability distribution for the next state.

Often notated with arcs between states

Markov Property

$$\begin{aligned} P(q_{t+1}=s_1 | q_t=s_2) &= 1/2 \\ P(q_{t+1}=s_2 | q_t=s_2) &= 1/2 \\ P(q_{t+1}=s_3 | q_t=s_2) &= 0 \end{aligned}$$

$$\begin{aligned} P(q_{t+1}=s_1 | q_t=s_1) &= 0 \\ P(q_{t+1}=s_2 | q_t=s_1) &= 0 \\ P(q_{t+1}=s_3 | q_t=s_1) &= 1 \end{aligned}$$



$N = 3$

$t=1$

$q_t=q_1=s_2$

q_{t+1} is conditionally independent of $\{ q_{t-1}, q_{t-2}, \dots, q_1, q_0 \}$ given q_t .

In other words:

$$P(q_{t+1} = s_j | q_t = s_i) =$$

$$P(q_{t+1} = s_j | q_t = s_i, \text{any earlier history})$$

$$\begin{aligned} P(q_{t+1}=s_1 | q_t=s_3) &= 1/3 \\ P(q_{t+1}=s_2 | q_t=s_3) &= 2/3 \\ P(q_{t+1}=s_3 | q_t=s_3) &= 0 \end{aligned}$$

Andrei Andreyevich Markov



Born	14 June 1856 N.S. Ryazan, Russian Empire
Died	20 July 1922 (aged 66) Petrograd, Russian SFSR
Residence	Russia
Nationality	Russian

Known for	Markov chains; Markov processes; stochastic processes
Scientific career	
Fields	Mathematician
Institutions	St. Petersburg University
Doctoral advisor	Pafnuty Chebyshev



Pafnuty Lvovich Chebyshev

Born	May 16, 1821 Akatovo, Kaluga Governorate, Russian Empire
Died	December 8, 1894 (aged 73) St. Petersburg, Russian Empire

As always, thanks to wikipedia

Markov Models

- Joint Distribution

$$\begin{aligned} p(\mathbf{X}) &= p(X_1, X_2, \dots, X_n) \\ &= p(X_1)p(X_2|X_1)p(X_3|X_2, X_1)\dots p(X_n|X_{n-1}, \dots, X_1) \\ &= \prod_{i=1}^n p(X_n|X_{n-1}, \dots, X_1) \end{aligned} \quad \text{Chain rule}$$

- Markov Assumption (m^{th} order)

$$p(\mathbf{X}) = \prod_{i=1}^n p(X_n|X_{n-1}, \dots, X_{n-m})$$

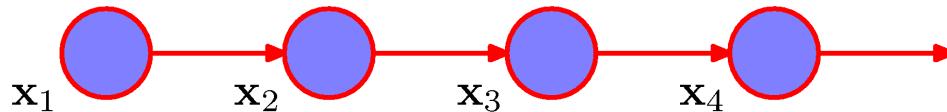
Current observation
only depends on past
 m observations

Markov Models

- Markov Assumption

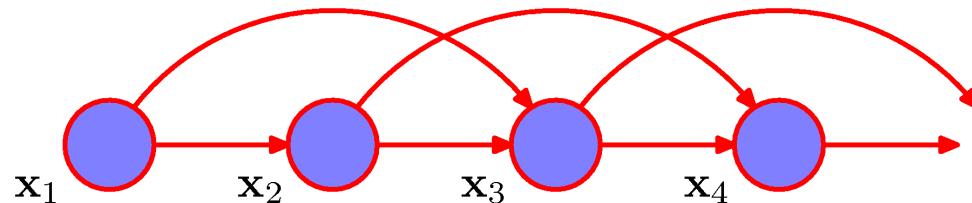
1st order

$$p(\mathbf{X}) = \prod_{i=1}^n p(X_n | X_{n-1})$$



2nd order

$$p(\mathbf{X}) = \prod_{i=1}^n p(X_n | X_{n-1}, X_{n-2})$$



Markov Models

- Markov Assumption

1st order

$$p(\mathbf{X}) = \prod_{i=1}^n p(X_n | X_{n-1})$$

parameters in
stationary model
K-ary variables

$O(K^2)$

mth order

$$p(\mathbf{X}) = \prod_{i=1}^n p(X_n | X_{n-1}, \dots, X_{n-m})$$

$O(K^{m+1})$

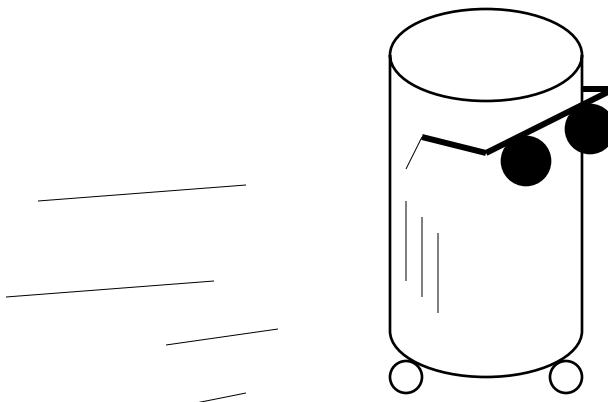
n-1th order

$$p(\mathbf{X}) = \prod_{i=1}^n p(X_n | X_{n-1}, \dots, X_1)$$

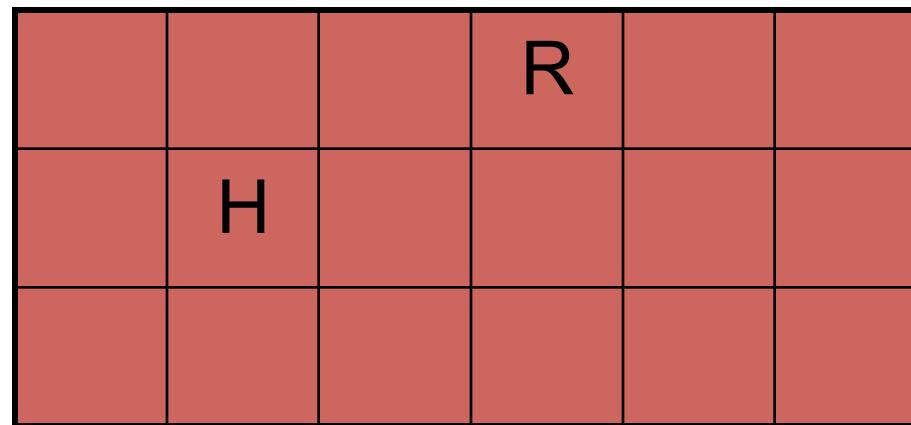
$O(K^n)$

≡ no assumptions – complete (but directed) graph

A Blind Robot



A Human and a Robot
wander around randomly
on a grid...



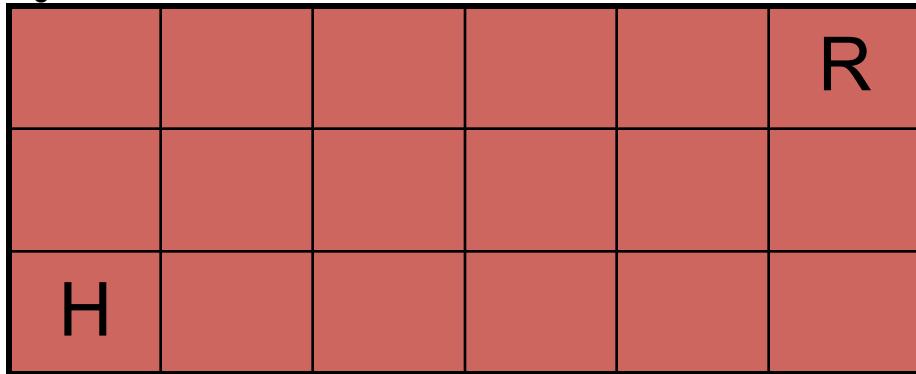
STATE $q =$

Location of Robot,
Location of Human

Note: N (num.
states) = $18 * 18 =$
324 (H and R)

Dynamics of System

$q_0 =$



Each timestep, the Human moves randomly to an adjacent cell. And Robot also moves randomly to an adjacent cell.

Typical Questions:

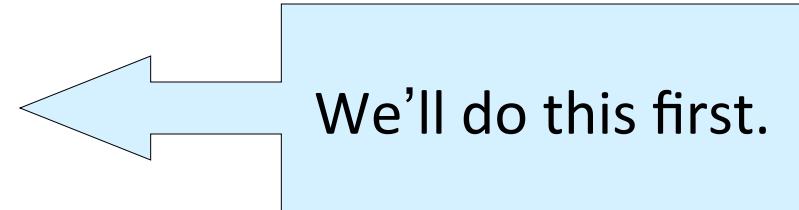
- “What is the expected time until the Human is met by the Robot?”
- “What is the probability that the Robot will hit the left wall before it meets the Human?”
- “What is the probability the Robot meets the Human on next time step?”

Example Question

“It’s currently time t , and human has not been met yet. What’s the probability of being met by the robot at time $t + 1$?”

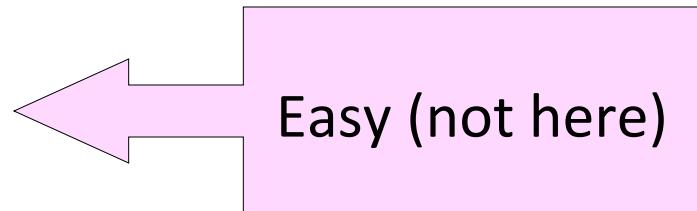
If robot is blind:

We can compute this in advance.



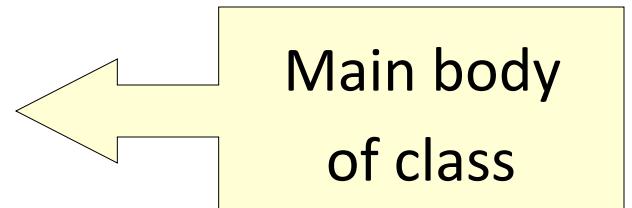
If robot is omnipotent:

(i.e., If robot knows state at time t),
can compute directly.



If robot has some sensors, but incomplete
state information ...

Hidden Markov Models are applicable!



What is $P(q_t = s)$? slow answer

Step 1: Work out how to compute $P(Q)$ for any path Q

$$= q_1 q_2 q_3 \dots q_t$$

Given we know the start state q_1 (i.e., $P(q_1) = 1$)

$$\begin{aligned} P(q_1 q_2 \dots q_t) &= P(q_1 q_2 \dots q_{t-1}) P(q_t | q_1 q_2 \dots q_{t-1}) \\ &= P(q_1 q_2 \dots q_{t-1}) P(q_t | q_{t-1}) \quad \text{Markov property} \\ &= P(q_2 | q_1) P(q_3 | q_2) \dots P(q_t | q_{t-1}) \end{aligned}$$

Step 2: Use this knowledge to get $P(q_t = s)$

$$P(q_t = s) = \sum_{Q \in \text{Paths of length } t \text{ that end in } s} P(Q)$$

Computation is exponential in t

What is $P(q_t = s)$? efficient answer

- For each state s_i , define

$$\begin{aligned} p_t(i) &= \text{Prob. state is } s_i \text{ at time } t \\ &= P(q_t = s_i) \end{aligned}$$

- Easy to do inductive definition

$$\forall i \quad p_0(i) =$$

$$\forall j \quad p_{t+1}(j) = P(q_{t+1} = s_j) =$$

What is $P(q_t = s)$? efficient answer

- For each state s_i , define

$$\begin{aligned} p_t(i) &= \text{Prob. state is } s_i \text{ at time } t \\ &= P(q_t = s_i) \end{aligned}$$

- Easy to do inductive definition

$$\forall i \quad p_0(i) = \begin{cases} 1 & \text{if } s_i \text{ is the start state} \\ 0 & \text{otherwise} \end{cases}$$

$$\forall j \quad p_{t+1}(j) = P(q_{t+1} = s_j) =$$

What is $P(q_t = s)$? efficient answer

- For each state s_i , define

$$\begin{aligned} p_t(i) &= \text{Prob. state is } s_i \text{ at time } t \\ &= P(q_t = s_i) \end{aligned}$$

- Easy to do inductive definition

$$\forall i \quad p_0(i) = \begin{cases} 1 & \text{if } s_i \text{ is the start state} \\ 0 & \text{otherwise} \end{cases}$$

$$\begin{aligned} \forall j \quad p_{t+1}(j) &= P(q_{t+1} = s_j) = \\ &\sum_{i=1}^N P(q_{t+1} = s_j \wedge q_t = s_i) = \end{aligned}$$

What is $P(q_t = s)$? efficient answer

- For each state s_i , define
 $p_t(i) = \text{Prob. state is } s_i \text{ at time } t$

$$= P(q_t = s_i)$$

- Easy to do inductive definition

$$\forall i \quad p_0(i) = \begin{cases} 1 & \text{if } s_i \text{ is the start state} \\ 0 & \text{otherwise} \end{cases}$$

$$\forall j \quad p_{t+1}(j) = P(q_{t+1} = s_j) = \sum_{i=1}^N P(q_{t+1} = s_j \wedge q_t = s_i) =$$

$$\sum_{i=1}^N P(q_{t+1} = s_j \mid q_t = s_i)P(q_t = s_i) = \sum_{i=1}^N a_{ij}p_t(i)$$

State Transition Probability

$$a_{ij} = P(q_{t+1} = s_j \mid q_t = s_i)$$

What is $P(q_t = s)$? efficient answer

- For each state s_i , define
 $p_t(i) = \text{Prob. state is } s_i \text{ at time } t$

$$= P(q_t = s_i)$$

- Easy to do inductive definition

$$\forall i \quad p_0(i) = \begin{cases} 1 & \text{if } s_i \text{ is the start state} \\ 0 & \text{otherwise} \end{cases}$$

$$\forall j \quad p_{t+1}(j) = P(q_{t+1} = s_j) =$$

$$\sum_{i=1}^N P(q_{t+1} = s_j \wedge q_t = s_i) =$$

$$\sum_{i=1}^N P(q_{t+1} = s_j \mid q_t = s_i) P(q_t = s_i) = \sum_{i=1}^N a_{ij} p_t(i)$$

- Computation is simple.
- Just fill in **this** table in **this** order:

t	$p_t(1)$	$p_t(2)$	\dots	$p_t(N)$
0	0	1		0
1				
:				
t_{final}				

What is $P(q_t = s)$? efficient answer

- For each state s_i , define

$$\begin{aligned} p_t(i) &= \text{Prob. state is } s_i \text{ at time } t \\ &= P(q_t = s_i) \end{aligned}$$

- Easy to do inductive definition

$$\forall i \quad p_0(i) = \begin{cases} 1 & \text{if } s_i \text{ is the start state} \\ 0 & \text{otherwise} \end{cases}$$

$$\forall j \quad p_{t+1}(j) = P(q_{t+1} = s_j) =$$

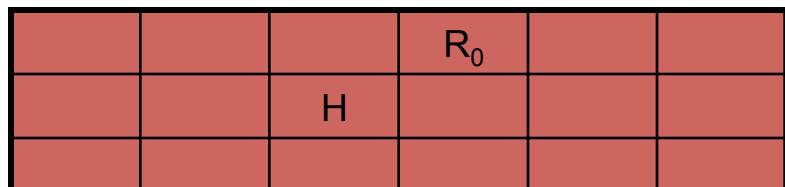
$$\sum_{i=1}^N P(q_{t+1} = s_j \wedge q_t = s_i) =$$

$$\sum_{i=1}^N P(q_{t+1} = s_j \mid q_t = s_i) P(q_t = s_i) = \sum_{i=1}^N a_{ij} p_t(i)$$

- Cost of computing $P_t(i)$ for all states S_i is now $O(t N^2)$
- The slow way was $O(N^t)$
- This was a simple example, as *Dynamic Programming*, used by HMMs algorithms.

Hidden State

- The previous example tried to estimate $P(q_t = s_i)$ unconditionally (using no observed evidence).
- Suppose we can **observe** something that is affected by the true state.
- Example: Proximity sensors (e.g., the contents of the 8 adjacent squares)



True state q_t



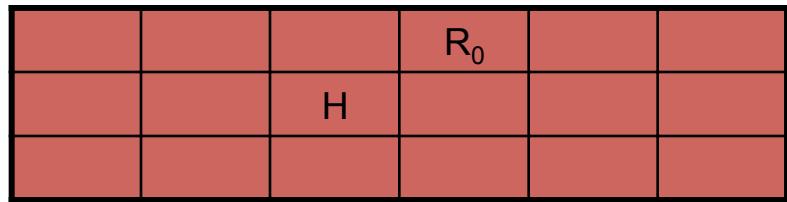
W	W	W
	®	
H		

W
denotes
“WALL”

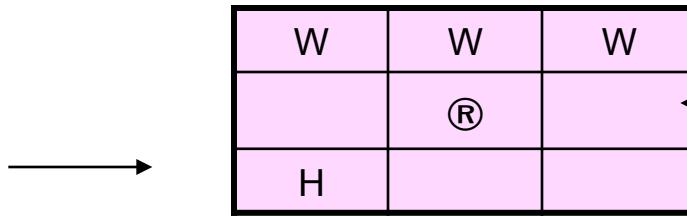
What the robot sees:
Observation O_t

Noisy Hidden State

- Realistic: Noisy Proximity sensors. (unreliable contents of the 8 adjacent squares)

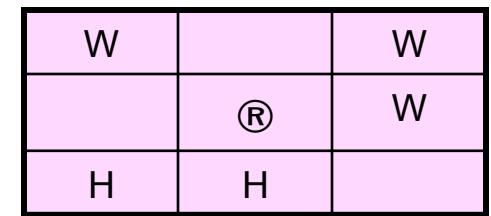


True state q_t



Uncorrupted Observation

W
denotes
“WALL”



What the robot sees:
Observation O_t

Noisy Hidden State

- Example: Noisy Proximity sensors. (unreliably tell us the contents of the 8 adjacent squares)

				R_0			2
			H				

True state q_t



W	W	W
	®	
H		

W
denotes
“WALL”

Uncorrupted Observation

O_t is noisily determined depending on the current state.

Assume that O_t is conditionally independent of $\{q_{t-1}, q_{t-2}, \dots, q_1, q_0, O_{t-1}, O_{t-2}, \dots, O_1, O_0\}$ given q_t .

i.e.:

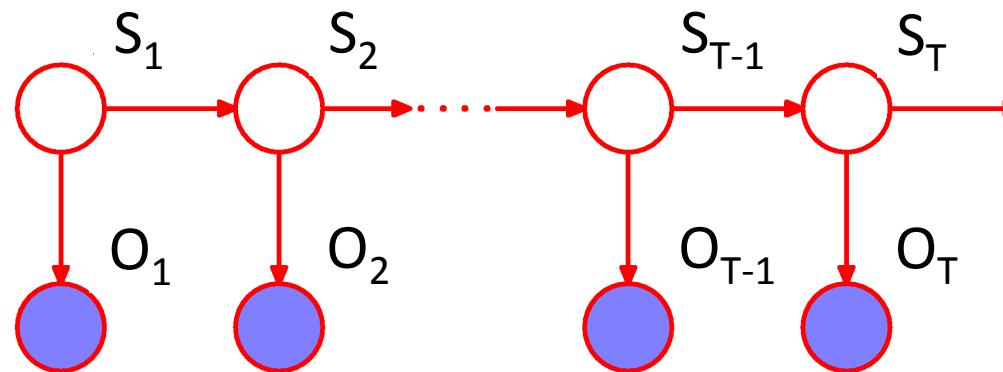
$$P(O_t = X | q_t = s_i) = P(O_t = X | q_t = s_i, \text{ any earlier history})$$

W		W
	®	
H	H	

What the robot sees:
Observation O_t

Hidden Markov Models

- Distributions that characterize sequential data with few parameters but are not limited by strong Markov assumptions.



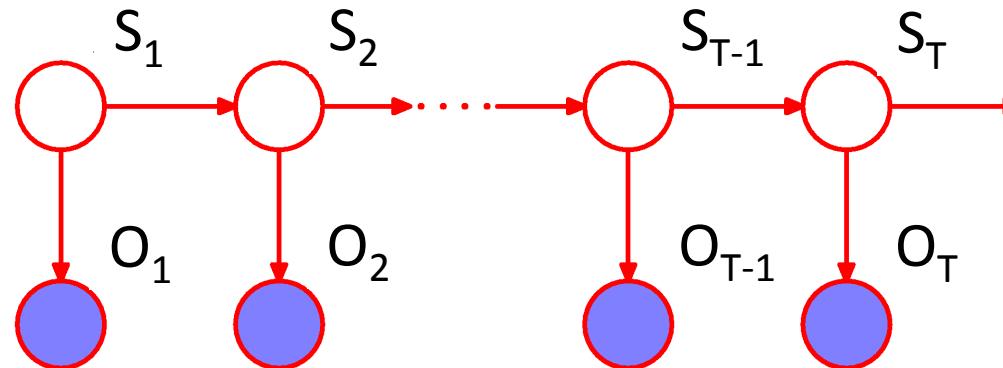
Observation space

$$O_t \in \{y_1, y_2, \dots, y_K\}$$

Hidden states

$$S_t \in \{1, \dots, I\}$$

Hidden Markov Models



$$p(S_1, \dots, S_T, O_1, \dots, O_T) = \prod_{t=1}^T p(O_t | S_t) \prod_{t=1}^T p(S_t | S_{t-1})$$

1st order Markov assumption on hidden states $\{S_t\}$ $t = 1, \dots, T$
(can be extended to higher order).

Hidden Markov Models

- Parameters – stationary/homogeneous Markov model (independent of time t)

Initial probabilities

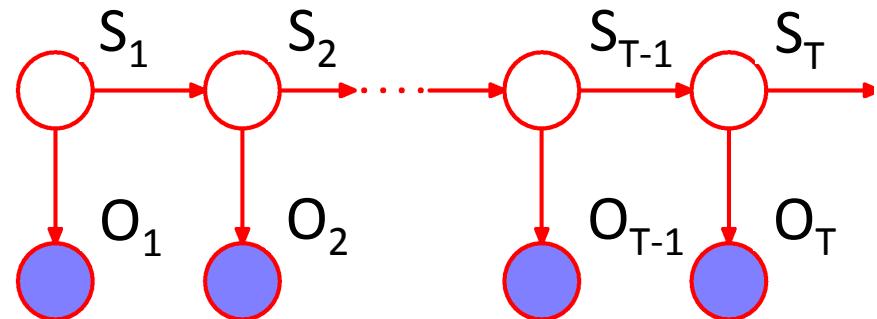
$$p(S_1 = i) = \pi_i$$

Transition probabilities

$$p(S_t = j | S_{t-1} = i) = p_{ij}$$

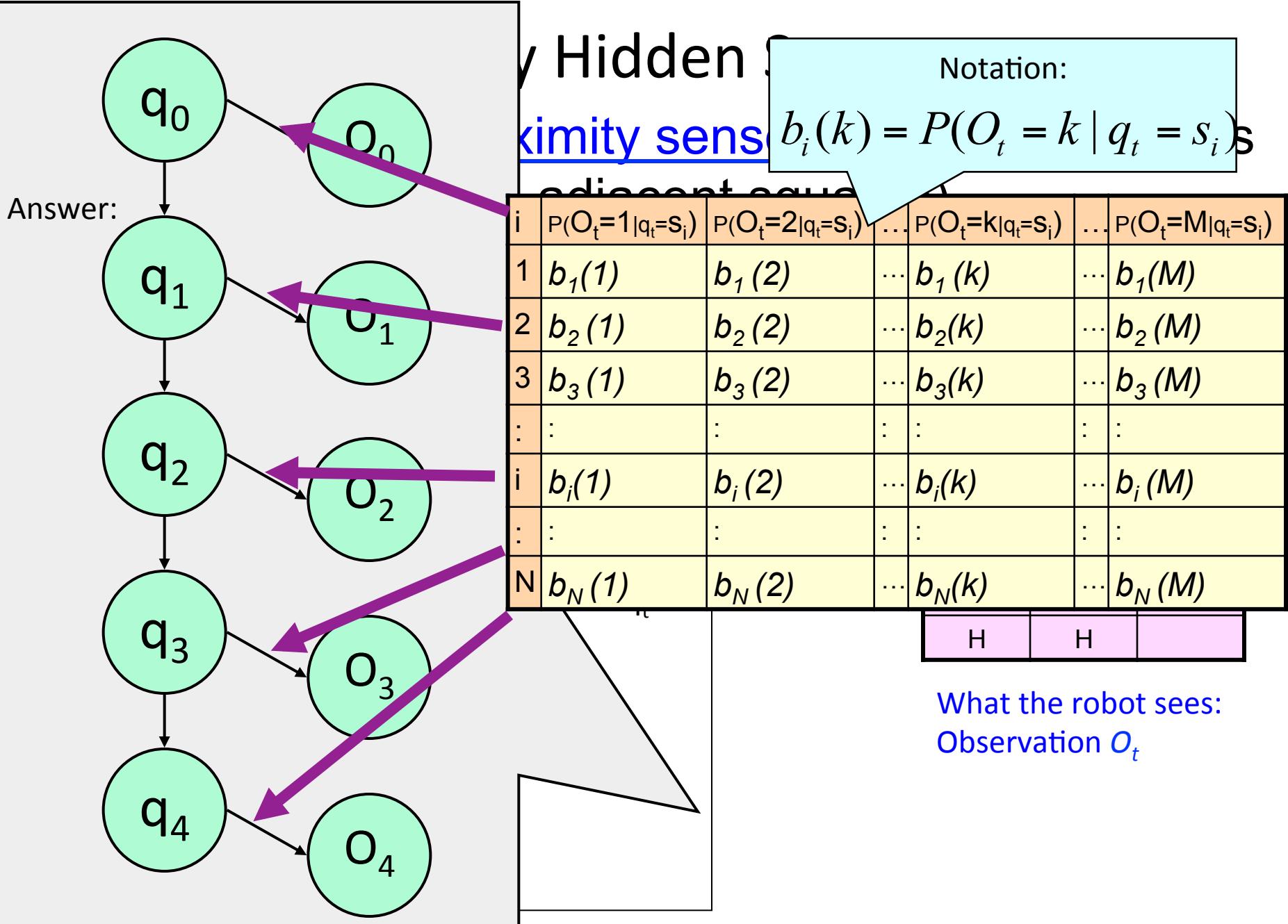
Emission probabilities

$$p(O_t = y | S_t = i) = q_i^y$$



$$p(\{S_t\}_{t=1}^T, \{O_t\}_{t=1}^T) =$$

$$p(S_1) \prod_{t=2}^T p(S_t | S_{t-1}) \prod_{t=1}^T p(O_t | S_t)$$



Hidden Markov Models

Our robot with noisy sensors is a good example of an HMM

- Question 1: State Estimation

What is $P(q_T = S_i \mid O_1 O_2 \dots O_T)$

It will turn out that a new cute D.P. trick will get this for us.

- Question 2: Most Probable Path

Given $O_1 O_2 \dots O_T$, what is the most probable path that I took?

And what is that probability?

Yet another famous D.P. trick, the VITERBI algorithm, gets this.

- Question 3: Learning HMMs:

Given $O_1 O_2 \dots O_T$, what is the maximum likelihood HMM that could have produced this string of observations?

Very very useful. Uses the E.M. Algorithm

Some Famous HMM Tasks

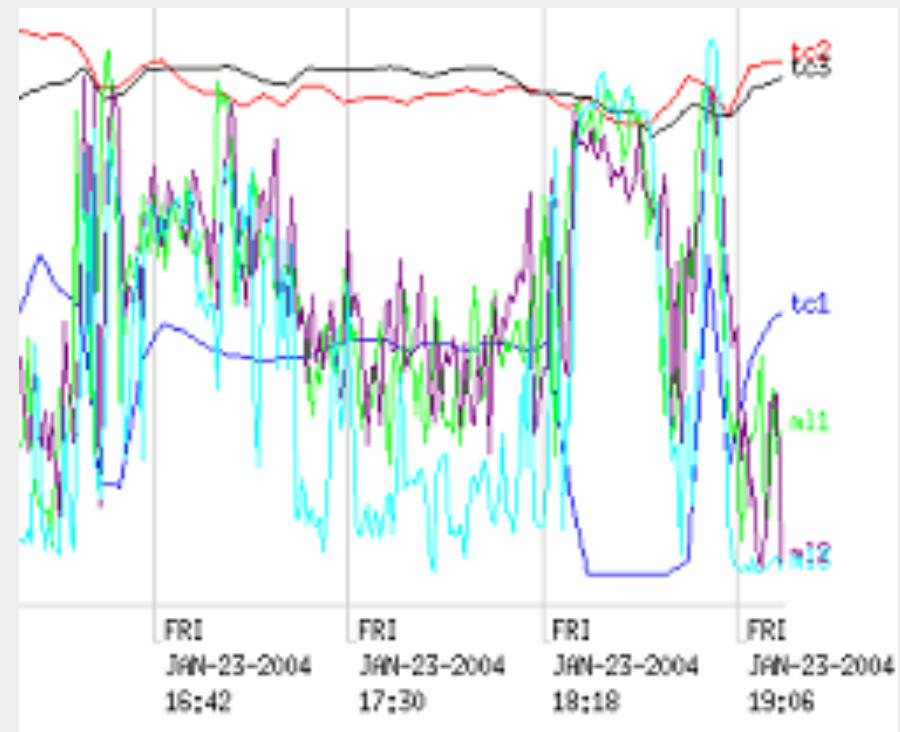
Question 1: State Estimation

What is $P(q_T = s_i | O_1 O_2 \dots O_t)$

Some Famous HMM Tasks

Question 1: State Estimation

What is $P(q_T = S_i | O_1, O_2, \dots, O_t)$

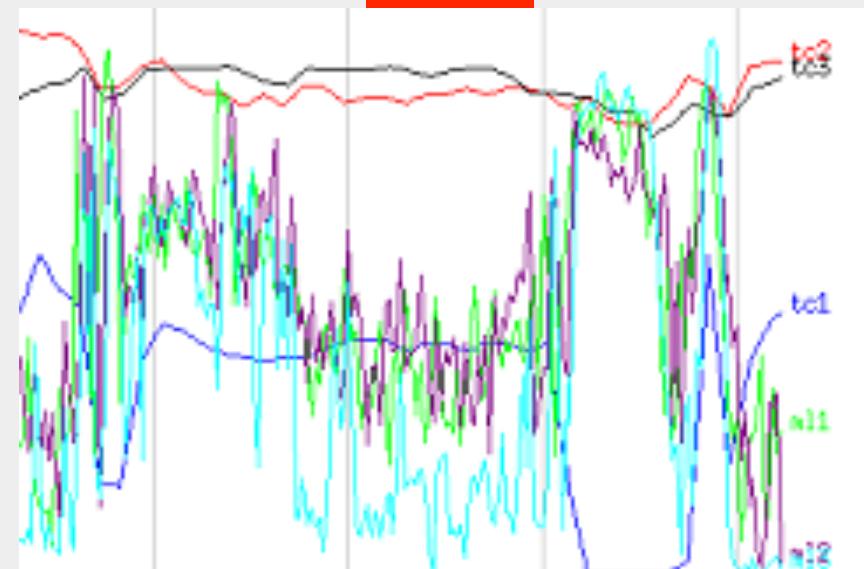
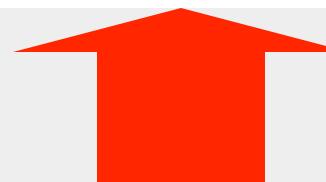
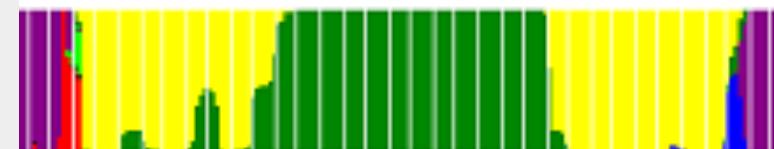


Some Famous HMM Tasks

Question 1: State Estimation

What is $P(q_T = S_i | O_1, O_2, \dots, O_t)$

bus class eat foo pc psw sleep stand swim walk



FRI
JAN-23-2004
16:42

FRI
JAN-23-2004
17:30

FRI
JAN-23-2004
18:18

FRI
JAN-23-2004
19:06

Some Famous HMM Tasks

Question 1: State Estimation

What is $P(q_T = s_i | O_1 O_2 \dots O_t)$

Question 2: Most Probable Path

Given $O_1 O_2 \dots O_T$, what is the most probable path that I took?

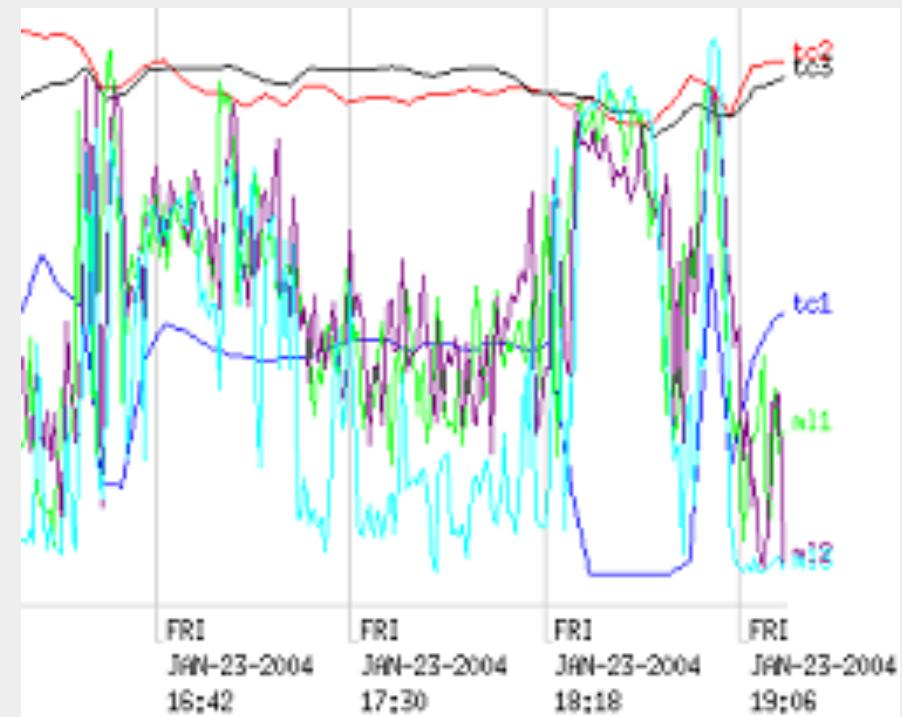
Some Famous HMM Tasks

Question 1: State Estimation

What is $P(q_T = S_i | O_1 O_2 \dots O_t)$

Question 2: Most Probable Path

Given $O_1 O_2 \dots O_T$, what is the most probable path that I took?



Some Famous HMM Tasks

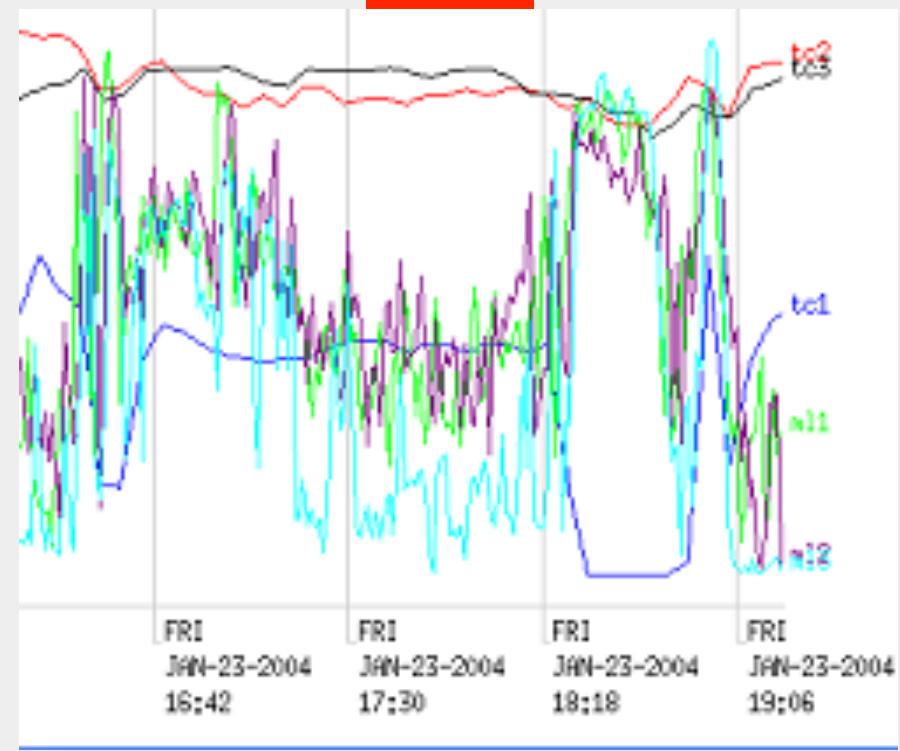
Question 1: State Estimation

What is $P(q_T=S_i | O_1O_2...O_t)$

Question 2: Most Probable Path

Given $O_1O_2...O_T$, what is the most probable path that I took?

Woke up at 8.35, Got on Bus at 9.46, Sat in lecture 10.05-11.22...



Some Famous HMM Tasks

Question 1: State Estimation

What is $P(q_T = s_i | O_1 O_2 \dots O_t)$

Question 2: Most Probable Path

Given $O_1 O_2 \dots O_T$, what is the most probable path that I took?

Question 3: Learning HMMs:

Given $O_1 O_2 \dots O_T$, what is the maximum likelihood HMM that could have produced this string of observations?

Some Famous Questions

Question 1: State Estimation

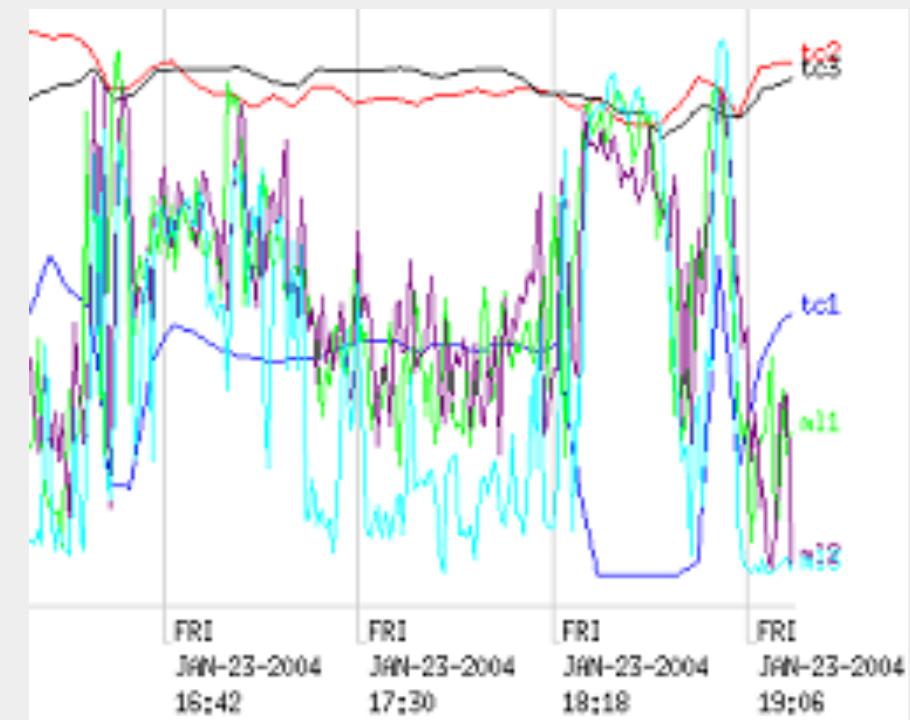
What is $P(q_T = S_i | O_1 O_2 \dots O_T)$

Question 2: Most Probable Path

Given $O_1 O_2 \dots O_T$, what is the most probable path that took?

Question 3: Learning HMMs:

Given $O_1 O_2 \dots O_T$, what is the maximum likelihood HMM that could have produced this string of observations?



Some Famous

Question 1: State Estimation

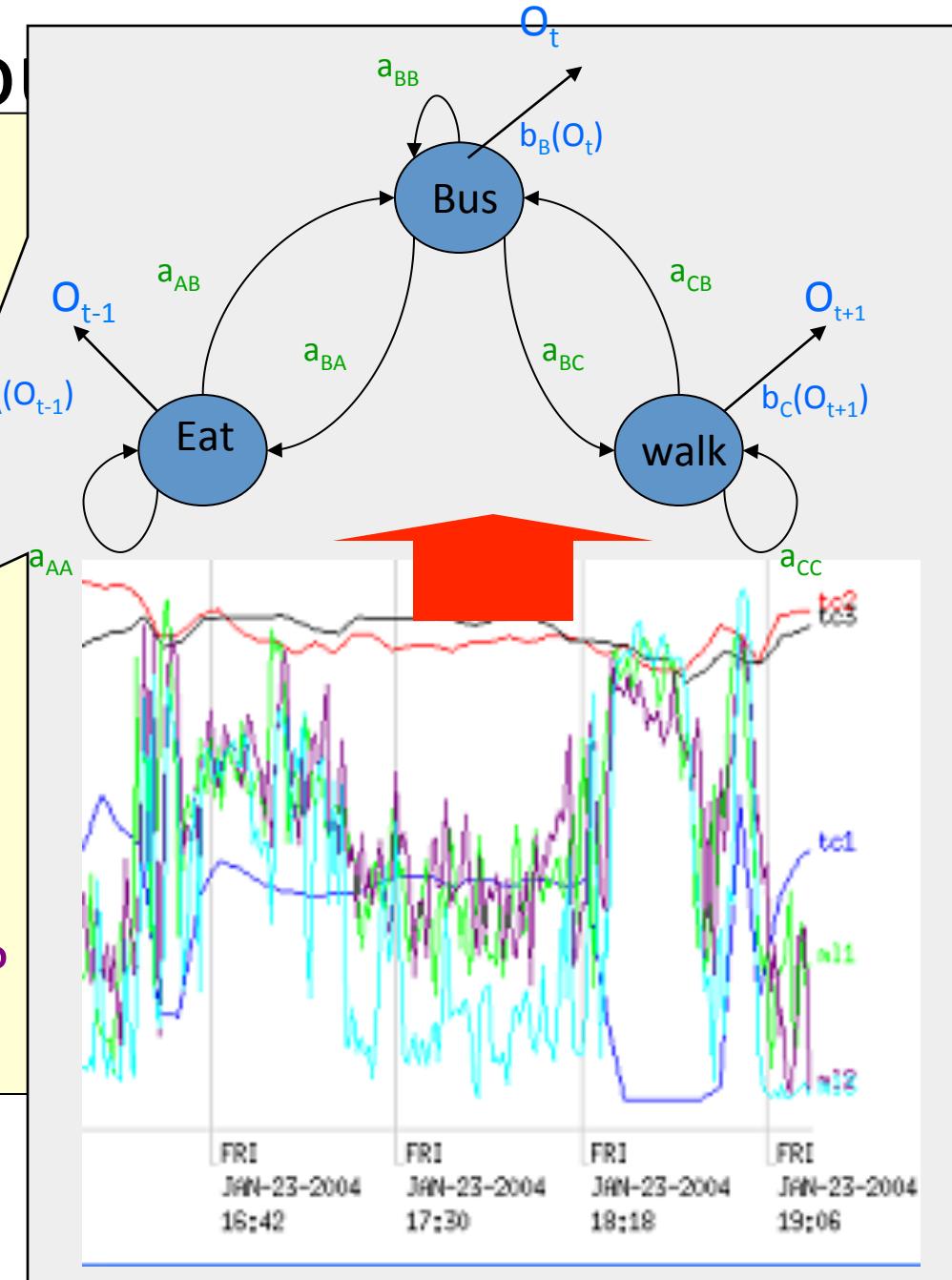
What is $P(q_T = S_i | O_1 O_2 \dots O_T)$

Question 2: Most Probable Path

Given $O_1 O_2 \dots O_T$, what is the most probable path that took?

Question 3: Learning HMMs:

Given $O_1 O_2 \dots O_T$, what is the maximum likelihood HMM that could have produced this string of observations?



HMM Notation (from Rabiner's Survey)

The states are labeled $S_1 S_2 \dots S_N$

For a particular trial....

Let T be the number of observations

T also the number of states passed through

$O = O_1 O_2 \dots O_T$ is the sequence of observations

$Q = q_1 q_2 \dots q_T$ is the notation for a path of states

$\lambda = \langle N, M, \{\pi_i\}, \{a_{ij}\}, \{b_i(j)\} \rangle$ is the specification of an HMM

*L. R. Rabiner, "A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition," Proc. of the IEEE, Vol.77, No.2, pp. 257--286, 1989.
Available from
<http://ieeexplore.ieee.org/iel5/5/698/00018626.pdf?arnumber=18626>

HMM Formal Definition

An HMM, λ , is a 5-tuple consisting of

- N the number of states
- M the number of possible observations
- $\{\pi_1, \pi_2, \dots, \pi_N\}$ The starting state probabilities

$$P(q_0 = S_i) = \pi_i$$

- $a_{11} \quad a_{22} \quad \dots \quad a_{1N}$
- $a_{21} \quad a_{22} \quad \dots \quad a_{2N}$
- $\vdots \quad \vdots \quad \vdots \quad \vdots$
- $a_{N1} \quad a_{N2} \quad \dots \quad a_{NN}$

The state transition probabilities

$$P(q_{t+1}=S_j | q_t=S_i) = a_{ij}$$

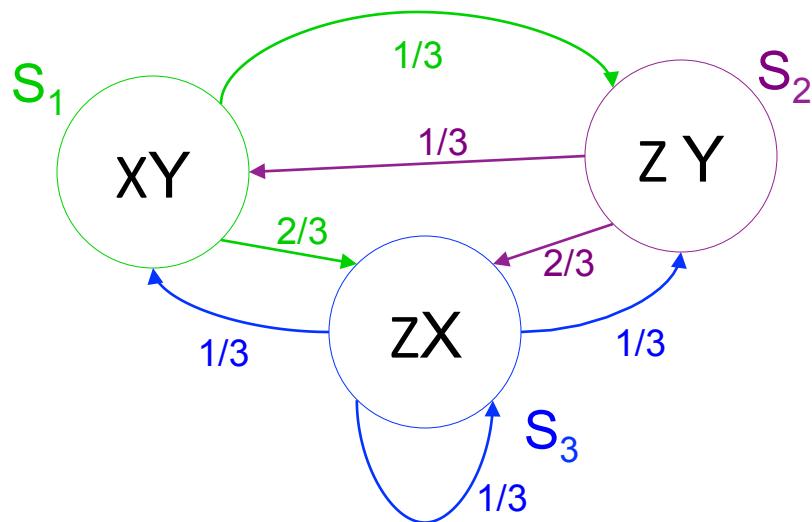
- $b_1(1) b_1(2) \quad \dots \quad b_1(M)$
- $b_2(1) b_2(2) \quad \dots \quad b_2(M)$
- $\vdots \quad \vdots \quad \vdots$
- $b_N(1) \quad b_N(2) \quad \dots \quad b_N(M)$

The observation probabilities

$$P(O_t=k | q_t=S_i) = b_i(k)$$

HMM Example

Start randomly in state 1 or 2
Choose one of the output symbols
in each state at random.



$$N = 3$$

$$M = 3$$

$$\pi_1 = \frac{1}{2}$$

$$\pi_2 = \frac{1}{2}$$

$$\pi_3 = 0$$

$$a_{11} = 0$$

$$a_{12} = \frac{1}{3}$$

$$a_{13} = \frac{2}{3}$$

$$a_{12} = \frac{1}{3}$$

$$a_{22} = 0$$

$$a_{13} = \frac{2}{3}$$

$$a_{13} = \frac{1}{3}$$

$$a_{32} = \frac{1}{3}$$

$$a_{13} = \frac{1}{3}$$

$$b_1(X) = \frac{1}{2} \quad b_1(Y) = \frac{1}{2}$$

$$b_2(X) = 0 \quad b_2(Y) = \frac{1}{2}$$

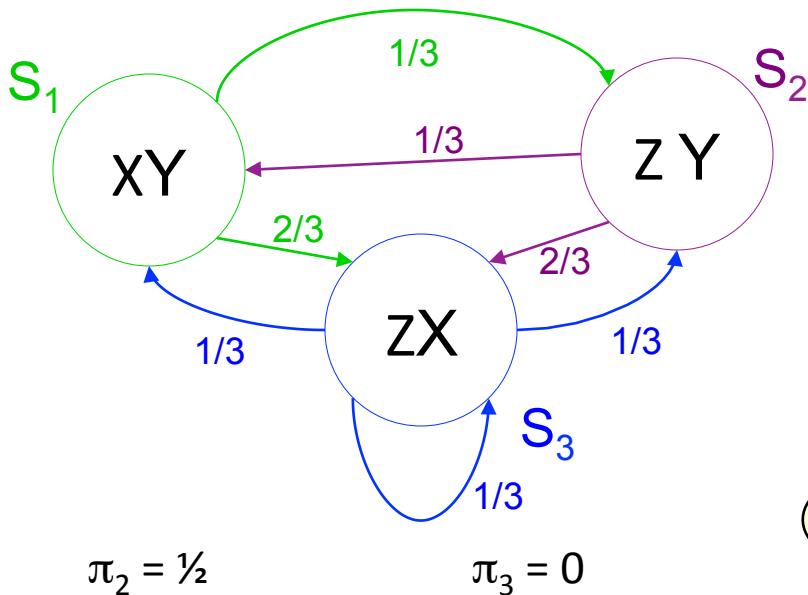
$$b_3(X) = \frac{1}{2} \quad b_3(Y) = 0$$

$$b_1(Z) = 0$$

$$b_2(Z) = \frac{1}{2}$$

$$b_3(Z) = \frac{1}{2}$$

HMM Example



$$N = 3$$

$$M = 3$$

$$\pi_1 = \frac{1}{2}$$

$$\pi_2 = \frac{1}{2}$$

$$\pi_3 = 0$$

$$a_{11} = 0$$

$$a_{12} = \frac{1}{3}$$

$$a_{13} = \frac{2}{3}$$

$$a_{12} = 0$$

$$a_{22} = 0$$

$$a_{32} = \frac{1}{3}$$

$$a_{13} = \frac{2}{3}$$

$$a_{13} = \frac{1}{3}$$

$$b_1(X) = \frac{1}{2}, b_1(Y) = \frac{1}{2}$$

$$b_2(X) = 0, b_2(Y) = \frac{1}{2}$$

$$b_3(X) = \frac{1}{2}, b_3(Y) = 0$$

$$b_1(Z) = 0$$

$$b_2(Z) = \frac{1}{2}$$

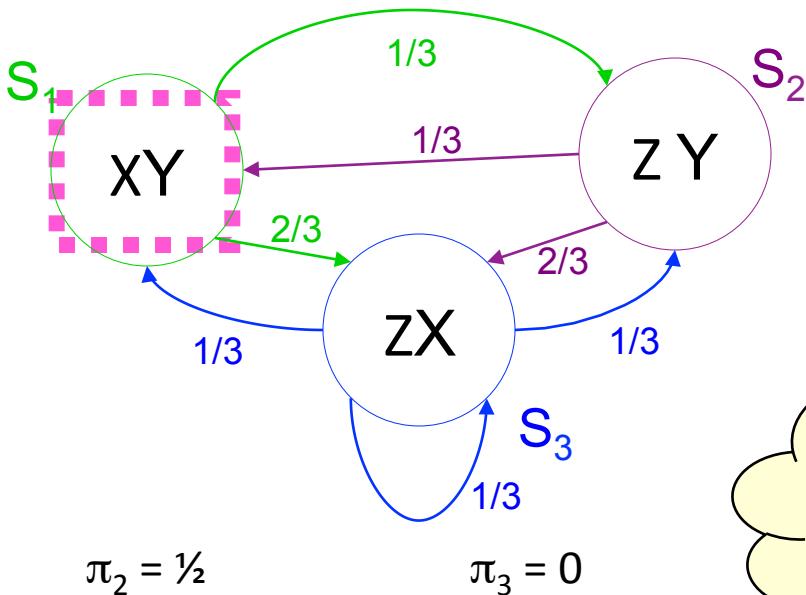
$$b_3(Z) = \frac{1}{2}$$

Start randomly in state 1 or 2
 Choose one of the output symbols
 in each state at random.
 Let's generate a sequence of
 observations:

50-50 choice
 between S_1 and S_2

$q_0 =$	○	$O_0 =$	—
$q_1 =$	—	$O_1 =$	—
$q_2 =$	—	$O_2 =$	—

HMM Example



$$N = 3$$

$$M = 3$$

$$\pi_1 = 1/2$$

$$\pi_2 = 1/2$$

$$\pi_3 = 0$$

$$a_{11} = 0$$

$$a_{12} = 1/3$$

$$a_{13} = 2/3$$

$$a_{12} = 1/3$$

$$a_{22} = 0$$

$$a_{13} = 2/3$$

$$a_{13} = 1/3$$

$$a_{32} = 1/3$$

$$a_{13} = 1/3$$

$$b_1(X) = 1/2 \quad b_1(Y) = 1/2$$

$$b_2(X) = 0 \quad b_2(Y) = 1/2$$

$$b_3(X) = 1/2 \quad b_3(Y) = 0$$

$$b_1(Z) = 0$$

$$b_2(Z) = 1/2$$

$$b_3(Z) = 1/2$$

Start randomly in state 1 or 2
 Choose one of the output symbols
 in each state at random.
 Let's generate a sequence of
 observations:

50-50 choice
 between X and Y

$q_0 =$	S_1	$O_0 =$	$\underline{\quad}$
$q_1 =$	$\underline{\quad}$	$O_1 =$	$\underline{\quad}$
$q_2 =$	$\underline{\quad}$	$O_2 =$	$\underline{\quad}$

HMM Example

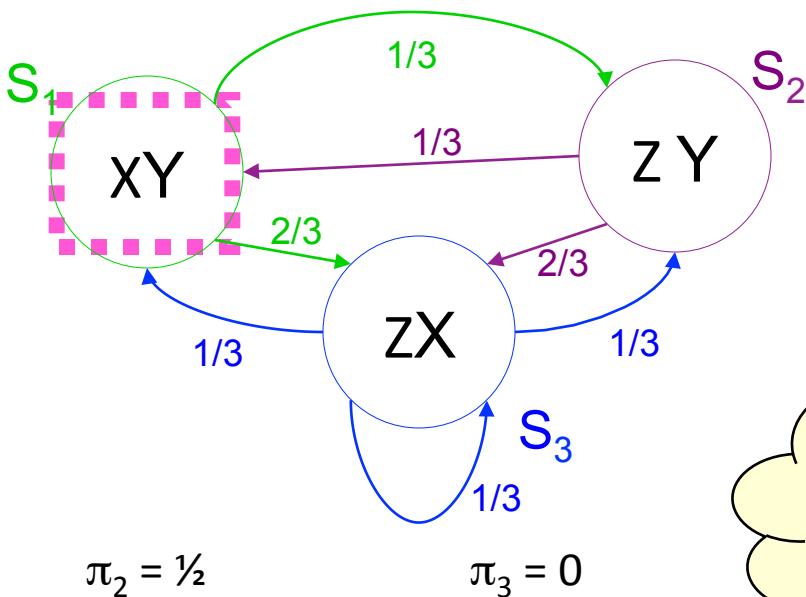
$$N = 3$$

$$M = 3$$

$$\pi_1 = \frac{1}{2}$$

$$\pi_2 = \frac{1}{2}$$

$$\pi_3 = 0$$



Start randomly in state 1 or 2
 Choose one of the output symbols
 in each state at random.
 Let's generate a sequence of
 observations:

Goto S_3 with
 probability $2/3$ or S_2
 with prob. $1/3$

$$a_{11} = 0$$

$$a_{12} = \frac{1}{3}$$

$$a_{13} = \frac{2}{3}$$

$$a_{12} = \frac{1}{3}$$

$$a_{22} = 0$$

$$a_{13} = \frac{2}{3}$$

$$a_{13} = \frac{1}{3}$$

$$a_{32} = \frac{1}{3}$$

$$a_{13} = \frac{1}{3}$$

$$b_1(X) = \frac{1}{2} \quad b_1(Y) = \frac{1}{2}$$

$$b_1(Z) = 0$$

$$b_2(X) = 0 \quad b_2(Y) = \frac{1}{2}$$

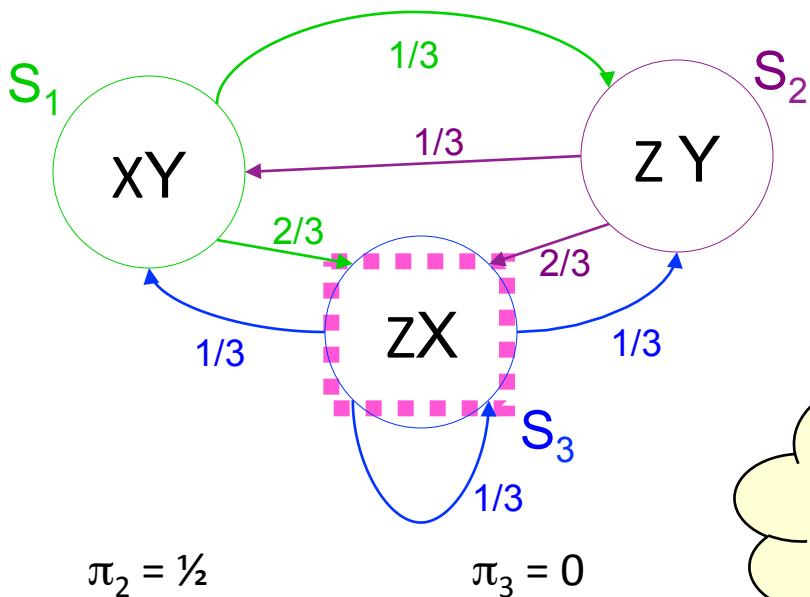
$$b_2(Z) = \frac{1}{2}$$

$$b_3(X) = \frac{1}{2} \quad b_3(Y) = 0$$

$$b_3(Z) = \frac{1}{2}$$

$q_0 =$	S_1	$O_0 =$	X
$q_1 =$	○	$O_1 =$	—
$q_2 =$	—	$O_2 =$	—

Here's an HMM



$a_{11} = 0$	$a_{12} = \frac{1}{3}$	$a_{13} = \frac{2}{3}$
$a_{12} = \frac{1}{3}$	$a_{22} = 0$	$a_{13} = \frac{2}{3}$
$a_{13} = \frac{1}{3}$	$a_{32} = \frac{1}{3}$	$a_{13} = \frac{1}{3}$

$b_1(X) = \frac{1}{2}$	$b_1(Y) = \frac{1}{2}$	$b_1(Z) = 0$
$b_2(X) = 0$	$b_2(Y) = \frac{1}{2}$	$b_2(Z) = \frac{1}{2}$
$b_3(X) = \frac{1}{2}$	$b_3(Y) = 0$	$b_3(Z) = \frac{1}{2}$

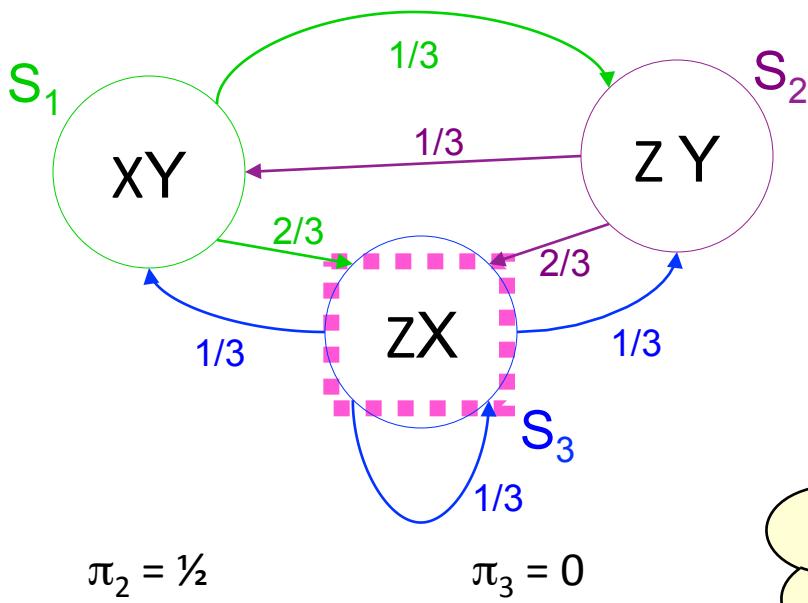
Start randomly in state 1 or 2
 Choose one of the output symbols
 in each state at random.
 Let's generate a sequence of
 observations:

50-50 choice
between Z and X



$q_0 =$	S_1	$O_0 =$	X
$q_1 =$	S_3	$O_1 =$	○
$q_2 =$	__	$O_2 =$	__

Here's an HMM



$$a_{11} = 0 \quad a_{12} = \frac{1}{3} \quad a_{13} = \frac{2}{3}$$

$$a_{12} = \frac{1}{3} \quad a_{22} = 0 \quad a_{13} = \frac{2}{3}$$

$$a_{13} = \frac{1}{3} \quad a_{32} = \frac{1}{3} \quad a_{13} = \frac{1}{3}$$

$$b_1(X) = \frac{1}{2} \quad b_1(Y) = \frac{1}{2}$$

$$b_2(X) = 0 \quad b_2(Y) = \frac{1}{2}$$

$$b_3(X) = \frac{1}{2} \quad b_3(Y) = 0$$

$$b_1(Z) = 0$$

$$b_2(Z) = \frac{1}{2}$$

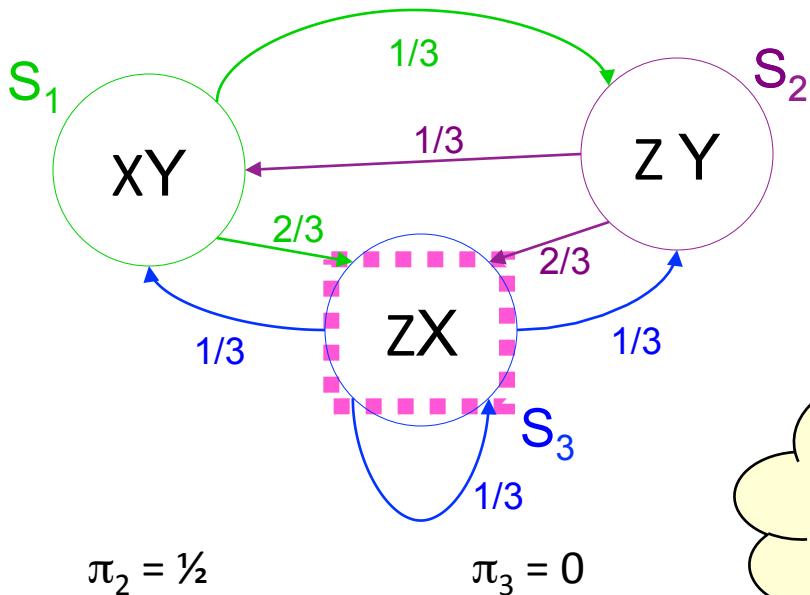
$$b_3(Z) = \frac{1}{2}$$

Start randomly in state 1 or 2
 Choose one of the output symbols
 in each state at random.
 Let's generate a sequence of
 observations:

Each of the three next states is equally likely

$q_0 =$	S_1	$O_0 =$	X
$q_1 =$	S_3	$O_1 =$	X
$q_2 =$	○	$O_2 =$	—

Here's an HMM



Start randomly in state 1 or 2
 Choose one of the output symbols
 in each state at random.
 Let's generate a sequence of
 observations:

50-50 choice
between Z and X

$$a_{11} = 0 \quad a_{12} = \frac{1}{3} \quad a_{13} = \frac{2}{3}$$

$$a_{12} = \frac{1}{3} \quad a_{22} = 0 \quad a_{13} = \frac{2}{3}$$

$$a_{13} = \frac{1}{3} \quad a_{32} = \frac{1}{3} \quad a_{13} = \frac{1}{3}$$

$$b_1(X) = \frac{1}{2} \quad b_1(Y) = \frac{1}{2}$$

$$b_2(X) = 0 \quad b_2(Y) = \frac{1}{2}$$

$$b_3(X) = \frac{1}{2} \quad b_3(Y) = 0$$

$$b_1(Z) = 0$$

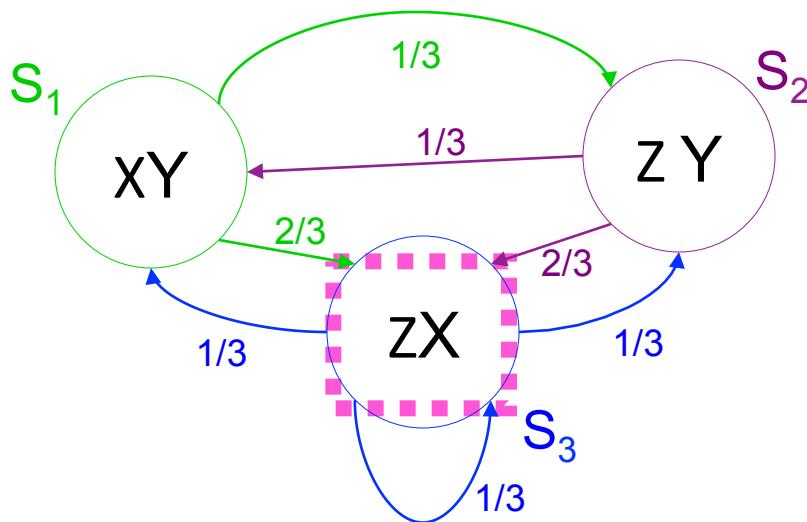
$$b_2(Z) = \frac{1}{2}$$

$$b_3(Z) = \frac{1}{2}$$

$q_0 =$	S_1	$O_0 =$	X
$q_1 =$	S_3	$O_1 =$	X
$q_2 =$	S_3	$O_2 =$	—

HMM Example

Start randomly in state 1 or 2
 Choose one of the output symbols
 in each state at random.
 Let's generate a sequence of
 observations:



$$N = 3$$

$$M = 3$$

$$\pi_1 = \frac{1}{2}$$

$$\pi_2 = \frac{1}{2}$$

$$\pi_3 = 0$$

$$a_{11} = 0$$

$$a_{12} = \frac{1}{3}$$

$$a_{13} = \frac{2}{3}$$

$$a_{12} = \frac{1}{3}$$

$$a_{22} = 0$$

$$a_{13} = \frac{2}{3}$$

$$a_{13} = \frac{1}{3}$$

$$a_{32} = \frac{1}{3}$$

$$a_{13} = \frac{1}{3}$$

$$b_1(X) = \frac{1}{2}, b_1(Y) = \frac{1}{2}$$

$$b_1(Z) = 0$$

$$b_2(X) = 0, b_2(Y) = \frac{1}{2}$$

$$b_2(Z) = \frac{1}{2}$$

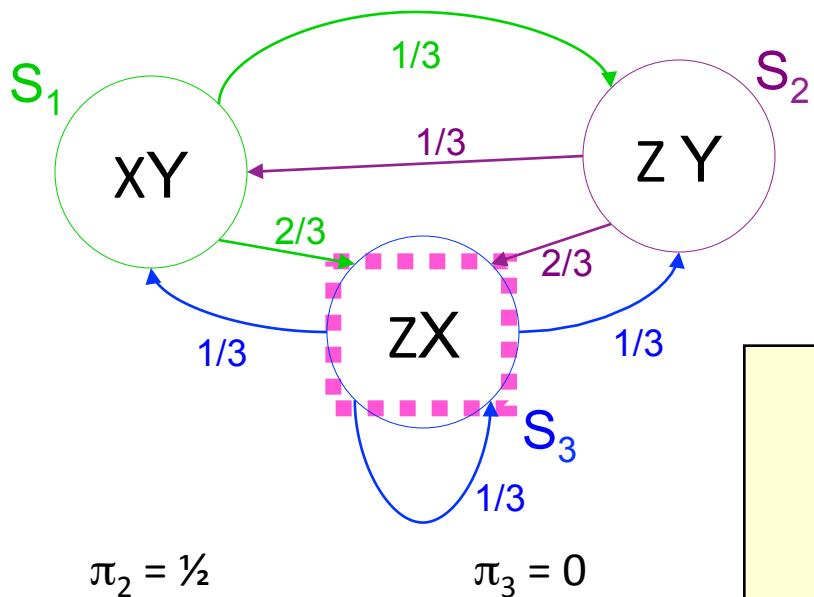
$$b_3(X) = \frac{1}{2}, b_3(Y) = 0$$

$$b_3(Z) = \frac{1}{2}$$

$q_0 =$	S_1	$O_0 =$	X
$q_1 =$	S_3	$O_1 =$	X
$q_2 =$	S_3	$O_2 =$	Z

State Estimation

Start randomly in state 1 or 2
 Choose one of the output symbols
 in each state at random.
 Let's generate a sequence of
 observations:



$$N = 3$$

$$M = 3$$

$$\pi_1 = \frac{1}{2}$$

$$a_{11} = 0$$

$$a_{12} = \frac{1}{3}$$

$$a_{13} = \frac{2}{3}$$

$$b_1(X) = \frac{1}{2} \quad b_1(Y) = \frac{1}{2}$$

$$b_2(X) = 0 \quad b_2(Y) = \frac{1}{2}$$

$$b_3(X) = \frac{1}{2} \quad b_3(Y) = 0$$

$$a_{12} = \frac{1}{3}$$

$$a_{22} = 0$$

$$a_{32} = \frac{1}{3}$$

$$a_{13} = \frac{2}{3}$$

$$a_{13} = \frac{2}{3}$$

$$a_{13} = \frac{1}{3}$$

$$b_1(Z) = 0$$

$$b_2(Z) = \frac{1}{2}$$

$$b_3(Z) = \frac{1}{2}$$

This is what the
 observer has to
 work with...

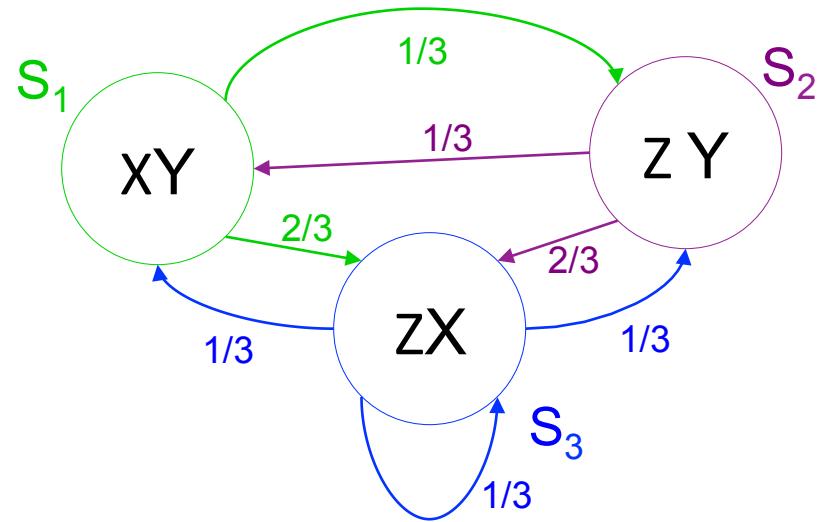
$q_0 =$?	$O_0 =$	X
$q_1 =$?	$O_1 =$	X
$q_2 =$?	$O_2 =$	Z

Prob. of a series of observations

What is $P(\mathbf{O}) = P(O_1 O_2 O_3) = P(O_1 = X \wedge O_2 = X \wedge O_3 = Z)$?

Slow way:

$$\begin{aligned} P(\mathbf{O}) &= \sum_{\mathbf{Q} \in \text{Paths of length 3}} P(\mathbf{O} \wedge \mathbf{Q}) \\ &= \sum_{\mathbf{Q} \in \text{Paths of length 3}} P(\mathbf{O} | \mathbf{Q})P(\mathbf{Q}) \end{aligned}$$



How do we compute $P(Q)$ for an arbitrary path Q ?

How do we compute $P(O|Q)$ for an arbitrary path Q ?

Prob. of a series of observations

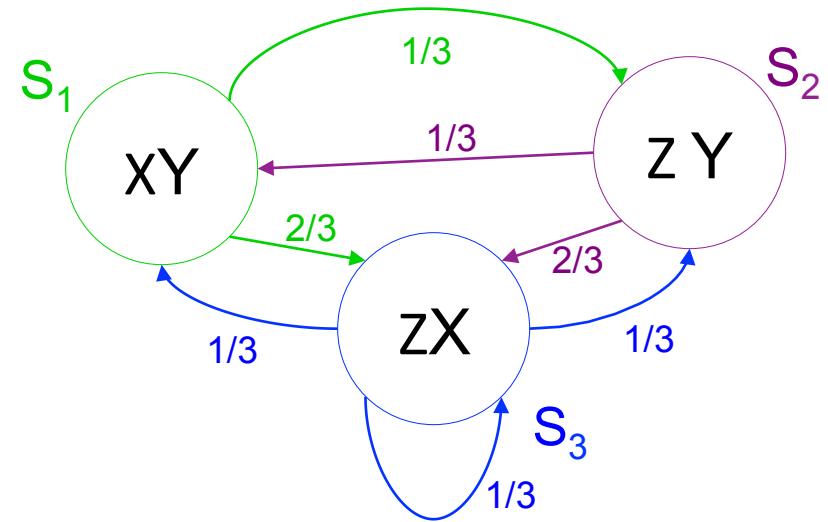
What is $P(\mathbf{O}) = P(O_1 O_2 O_3) = P(O_1 = X \wedge O_2 = Y \wedge O_3 = Z)$?

Slow way:

$$\begin{aligned} P(\mathbf{O}) &= \sum_{Q \in \text{Paths of length 3}} P(\mathbf{O} \wedge Q) \\ &= \sum_{Q \in \text{Paths of length 3}} P(\mathbf{O} | Q)P(Q) \end{aligned}$$

How do we compute $P(Q)$ for an arbitrary path Q ?

How do we compute $P(\mathbf{O}|Q)$ for an arbitrary path Q ?



$$P(Q) = P(q_1, q_2, q_3)$$

$$= P(q_1) P(q_2, q_3 | q_1) \text{ (chain rule)}$$

$$= P(q_1) P(q_2 | q_1) P(q_3 | q_2, q_1) \text{ (chain)}$$

$$= P(q_1) P(q_2 | q_1) P(q_3 | q_2) \text{ (why?)}$$

Example in the case $Q = S_1 S_3 S_3$:

$$= 1/2 * 2/3 * 1/3 = 1/9$$

Probability of a series of observations

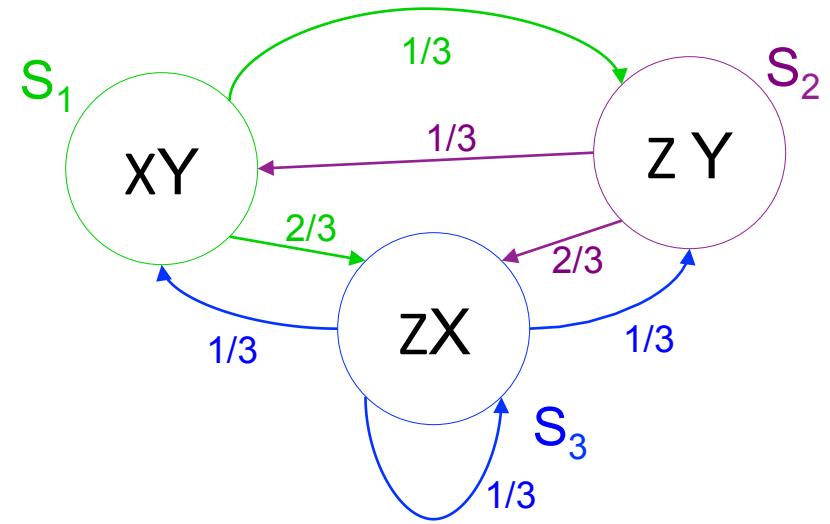
What is $P(\mathbf{O}) = P(O_1 O_2 O_3) = P(O_1 = X \wedge O_2 = Y \wedge O_3 = Z)$?

Slow way:

$$\begin{aligned} P(\mathbf{O}) &= \sum_{Q \in \text{Paths of length 3}} P(\mathbf{O} \wedge Q) \\ &= \sum_{Q \in \text{Paths of length 3}} P(\mathbf{O} | Q)P(Q) \end{aligned}$$

How do we compute $P(Q)$ for an arbitrary path Q ?

How do we compute $P(\mathbf{O}|Q)$ for an arbitrary path Q ?



$$\begin{aligned} P(\mathbf{O} | Q) &= P(O_1 O_2 O_3 | q_1 q_2 q_3) \\ &= P(O_1 | q_1) P(O_2 | q_2) P(O_3 | q_3) \end{aligned}$$

Example in the case $Q = S_1 S_3 S_3$:

$$\begin{aligned} &= P(X | S_1) P(Y | S_3) P(Z | S_3) = \\ &= 1/2 * 1/2 * 1/2 = 1/8 \end{aligned}$$

Probability of a series of observations

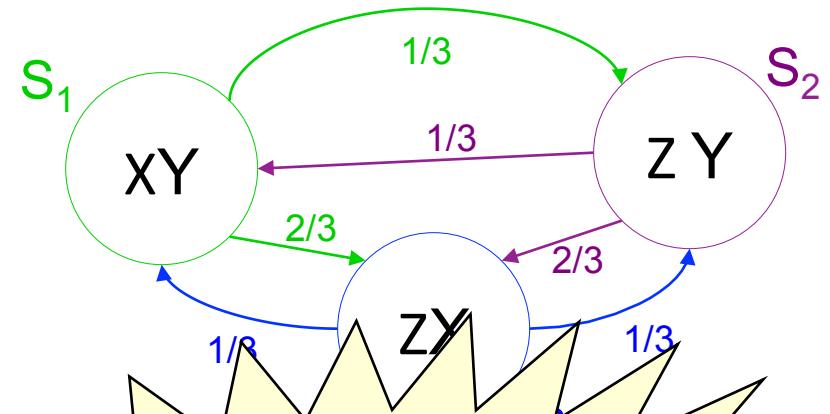
What is $P(\mathbf{O}) = P(O_1 O_2 O_3) = P(O_1 = X \wedge O_2 = X \wedge O_3 = Z)$?

Slow way:

$$\begin{aligned} P(\mathbf{O}) &= \sum_{Q \in \text{Paths of length 3}} P(\mathbf{O} \wedge Q) \\ &= \sum_{Q \in \text{Paths of length 3}} P(\mathbf{O} | Q)P(Q) \end{aligned}$$

How do we compute $P(Q)$ for an arbitrary path Q ?

How do we compute $P(\mathbf{O}|Q)$ for an arbitrary path Q ?



$P(\mathbf{O})$ would need 27 $P(Q)$ computations and 27 $P(\mathbf{O}|Q)$ computations

A sequence of 20 observations would need $3^{20} = 3.5$ billion computations and 3.5 billion $P(\mathbf{O}|Q)$ computations

So let's be smarter...

The Prob. of a given series of observations, non-exponential-cost-style

Given observations $O_1 O_2 \dots O_T$

Define

$$\alpha_t(i) = P(O_1 O_2 \dots O_t \wedge q_t = S_i | \lambda) \quad \text{where } 1 \leq t \leq T$$

$\alpha_t(i)$ = Probability that, in a random trial, we would have

- seen the first t observations
- and ended up in S_i as the state visited at time t

$\alpha_t(i)$: easy to define recursively

$$\alpha_t(i) = P(O_1 O_2 \dots O_T \wedge q_t = S_i | \lambda)$$

$$\begin{aligned}\alpha_1(i) &= P(O_1 \wedge q_1 = S_i) \\ &= P(q_1 = S_i) P(O_1 | q_1 = S_i) \\ &= b_i(O_1) \pi_i\end{aligned}$$

$$\begin{aligned}\alpha_{t+1}(j) &= P(O_1 O_2 \dots O_t O_{t+1} \wedge q_{t+1} = S_j) \\ &= \dots\end{aligned}$$

$\alpha_t(i)$: easy to define recursively

$$\alpha_t(i) = P(O_1 O_2 \dots O_T \wedge q_t = S_i | \lambda)$$

$$\begin{aligned}\alpha_1(i) &= P(O_1 \wedge q_1 = S_i) \\ &= P(q_1 = S_i) P(O_1 | q_1 = S_i) \\ &= b_i(O_1) \pi_i\end{aligned}$$

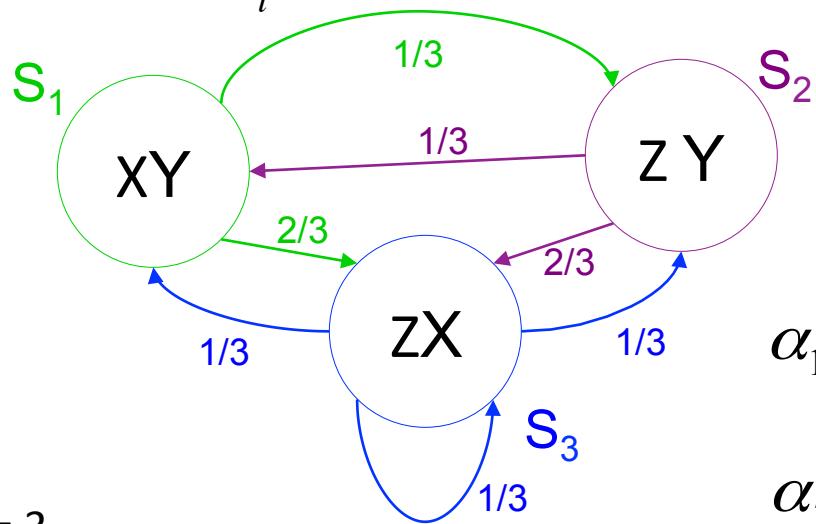
$$\begin{aligned}\alpha_{t+1}(j) &= P(O_1 O_2 \dots O_t O_{t+1} \wedge q_{t+1} = S_j) \\ &= \sum_{i=1}^N P(O_1 O_2 \dots O_t \wedge q_t = S_i \wedge O_{t+1} \wedge q_{t+1} = S_j) \\ &= \sum_{i=1}^N P(O_{t+1}, q_{t+1} = S_j | O_1 O_2 \dots O_t \wedge q_t = S_i) P(O_1 O_2 \dots O_t \wedge q_t = S_i) \\ &= \sum_i P(O_{t+1}, q_{t+1} = S_j | q_t = S_i) \alpha_t(i) \\ &= \sum_i P(q_{t+1} = S_j | q_t = S_i) P(O_{t+1} | q_{t+1} = S_j) \alpha_t(i) \\ &= \sum_i a_{ij} b_j(O_{t+1}) \alpha_t(i)\end{aligned}$$

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

$$\alpha_1(i) = b_i(O_1) \pi_i$$

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

Assume $O_1 O_2 O_3 = X X Z$



$$N = 3$$

$$M = 3$$

$$\pi_1 = \frac{1}{2} \quad \pi_2 = \frac{1}{2} \quad \pi_3 = 0$$

$$\alpha_1(1) = \frac{1}{4}$$

$$\alpha_2(1) = 0$$

$$\alpha_3(1) = 0$$

$$\alpha_1(2) = 0$$

$$\alpha_2(2) = 0$$

$$\alpha_3(2) = \frac{1}{72}$$

$$\alpha_1(3) = 0$$

$$\alpha_2(3) = \frac{1}{12}$$

$$\alpha_3(3) = \frac{1}{72}$$

$b_1(X) = \frac{1}{2}$	$b_1(Y) = \frac{1}{2}$	$b_1(Z) = 0$
$b_2(X) = 0$	$b_2(Y) = \frac{1}{2}$	$b_2(Z) = \frac{1}{2}$
$b_3(X) = \frac{1}{2}$	$b_3(Y) = 0$	$b_3(Z) = \frac{1}{2}$

Easy Question

We can cheaply compute

$$a_t(i) = P(O_1 O_2 \dots O_t \wedge q_t = S_i)$$

(How) can we cheaply compute

$$P(O_1 O_2 \dots O_t) ?$$

(How) can we cheaply compute

$$P(q_t = S_i | O_1 O_2 \dots O_t)$$

Following Easy Questions

We can cheaply compute

$$\alpha_t(i) = P(O_1 O_2 \dots O_t \wedge q_t = S_i)$$

(How) can we cheaply compute

$$P(O_1 O_2 \dots O_t) \quad ?$$

$$\sum_{i=1}^N \alpha_t(i)$$

(How) can we cheaply compute

$$P(q_t = S_i | O_1 O_2 \dots O_t)$$

$$\frac{\alpha_t(i)}{\sum_{j=1}^N \alpha_t(j)}$$

New question: Most probable **path** given observations

What is the most probable path given $O_1O_2\dots O_T$, i.e.

What is $\underset{Q}{\operatorname{argmax}} P(Q|O_1O_2\dots O_T)$?

Slow answer:

$$\underset{Q}{\operatorname{argmax}} P(Q|O_1O_2\dots O_T)$$

$$= \underset{Q}{\operatorname{argmax}} \frac{P(O_1O_2\dots O_T|Q)P(Q)}{P(O_1O_2\dots O_T)}$$

$$= \underset{Q}{\operatorname{argmax}} P(O_1O_2\dots O_T|Q)P(Q)$$

Efficient MPP computation

We compute the following variables:

$$\delta_t(i) = \max_{q_1 q_2 \dots q_{t-1}} P(q_1 q_2 \dots q_{t-1} \wedge q_t = S_i \wedge O_1 \dots O_t)$$

= The probability of the PATH of length t-1 with the maximum chance of doing all these things:

...the path ITSELF OCCURING

and

...the path ENDING UP IN STATE S_i

and

...the path PRODUCING OUTPUT $O_1 \dots O_t$

DEFINE: $mpp_t(i)$ = that most probable path

So: $\delta_t(i) = \text{Prob}(mpp_t(i))$

Andrew J. Viterbi^[1]



Born	Andrea Giacomo Viterbi March 9, 1935 (age 83) Bergamo, Italy
Citizenship	Italian, American
Education	Massachusetts Institute of Technology (BS, MS) University of Southern California (PhD)

The Viterbi Algorithm

$$\delta_t(i) = \max_{q_1 q_2 \dots q_{t-1}} P(q_1 q_2 \dots q_{t-1} \wedge q_t = S_i \wedge O_1 O_2 \dots O_t)$$

$$mpp_t(i) = \arg \max_{q_1 q_2 \dots q_{t-1}} P(q_1 q_2 \dots q_{t-1} \wedge q_t = S_i \wedge O_1 O_2 \dots O_t)$$

$$\delta_1(i) = \max_{\text{one choice}} P(q_1 = S_i \wedge O_1)$$

$$= P(q_1 = S_i) P(O_1 | q_1 = S_i)$$

$$= \pi_i b_i(O_1)$$

The Viterbi Algorithm

$$\delta_t(i) = \max_{q_1 q_2 \dots q_{t-1}} P(q_1 q_2 \dots q_{t-1} \wedge q_t = S_i \wedge O_1 O_2 \dots O_t)$$

$$mpp_t(i) = \arg \max_{q_1 q_2 \dots q_{t-1}} P(q_1 q_2 \dots q_{t-1} \wedge q_t = S_i \wedge O_1 O_2 \dots O_t)$$

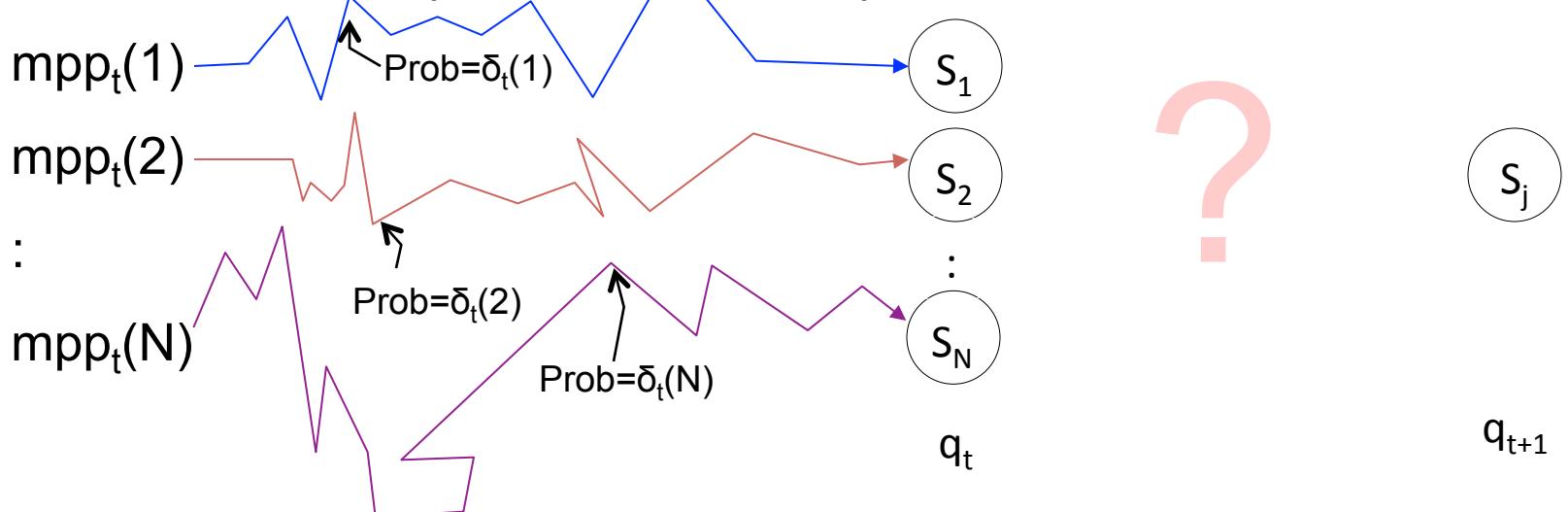
$$\delta_1(i) = \max_{\text{one choice}} P(q_1 = S_i \wedge O_1)$$

$$= P(q_1 = S_i) P(O_1 | q_1 = S_i)$$

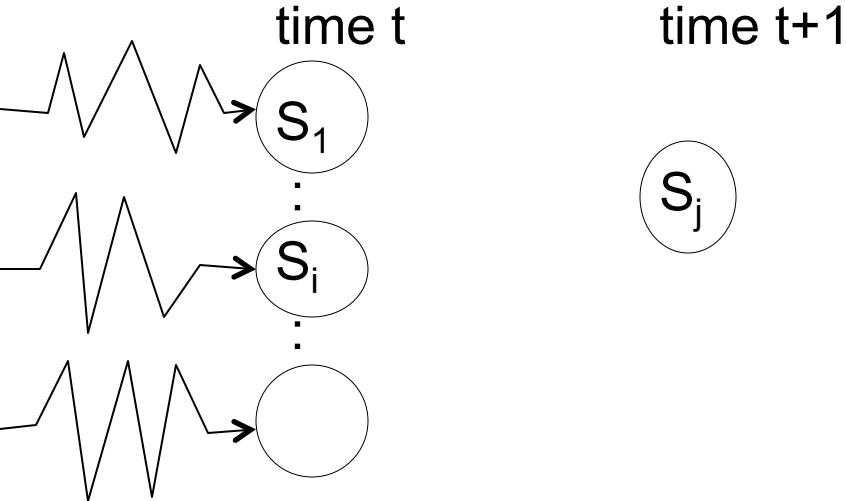
$$= \pi_i b_i(O_1)$$

If we have all the $\delta_t(i)$'s and $mpp_t(i)$'s for all i.

HOW TO GET $\delta_{t+1}(j)$ and $mpp_{t+1}(j)$?

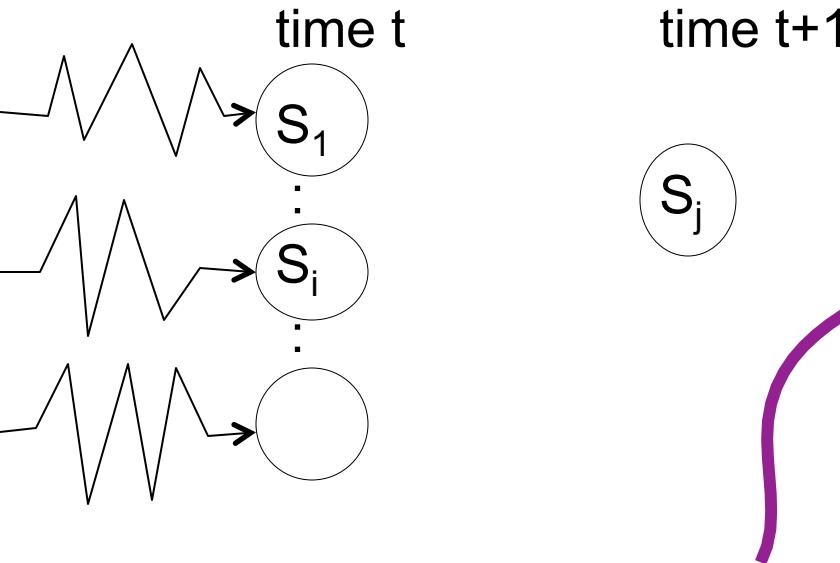


The Viterbi Algorithm



The most prob path with last two states $S_i \ S_j$ is the most prob path to S_i , followed by transition $S_i \rightarrow S_j$

The Viterbi Algorithm



The most prob path with last two states $S_i \ S_j$
is
the most prob path to S_i ,
followed by transition $S_i \rightarrow S_j$

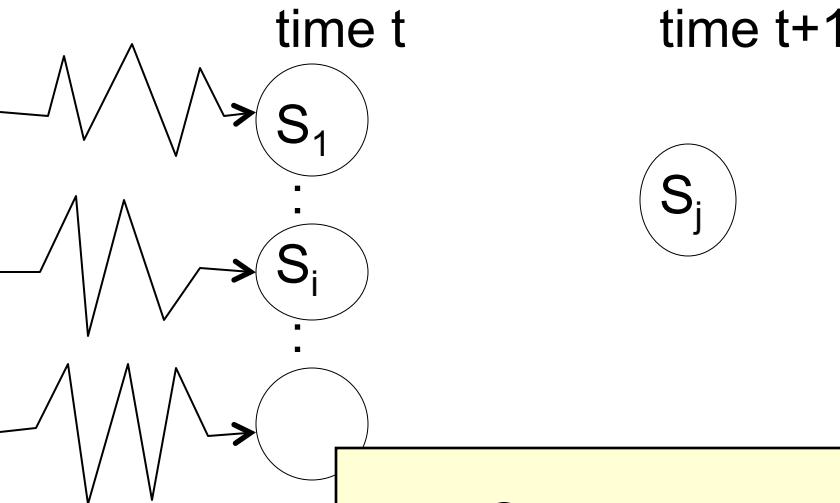
What is the prob of that path?

$$\begin{aligned} & \delta_t(i) \times P(S_i \rightarrow S_j \wedge O_{t+1} | \lambda) \\ &= \delta_t(i) a_{ij} b_j(O_{t+1}) \end{aligned}$$

SO The most probable path to S_j has S_{i^*} as its penultimate state

where $i^* = \operatorname{argmax}_i \delta_t(i) a_{ij} b_j(O_{t+1})$

The Viterbi Algorithm



Summary:

$$\delta_{t+1}(j) = \delta_t(i^*) a_{ij} b_j(O_{t+1})$$

$$mpp_{t+1}(j) = mpp_{t+1}(i^*) S_{i^*}$$

where $i^* = \text{argmax } \delta_t(i) a_{ij} b_j(O_{t+1})$

What is Viterbi Algorithm used for?

Classic Example

Speech recognition:

Signal → words

HMM → observable is signal

→ Hidden state is part of word
formation

What is the most probable word given this signal?

UTTERLY GROSS SIMPLIFICATION

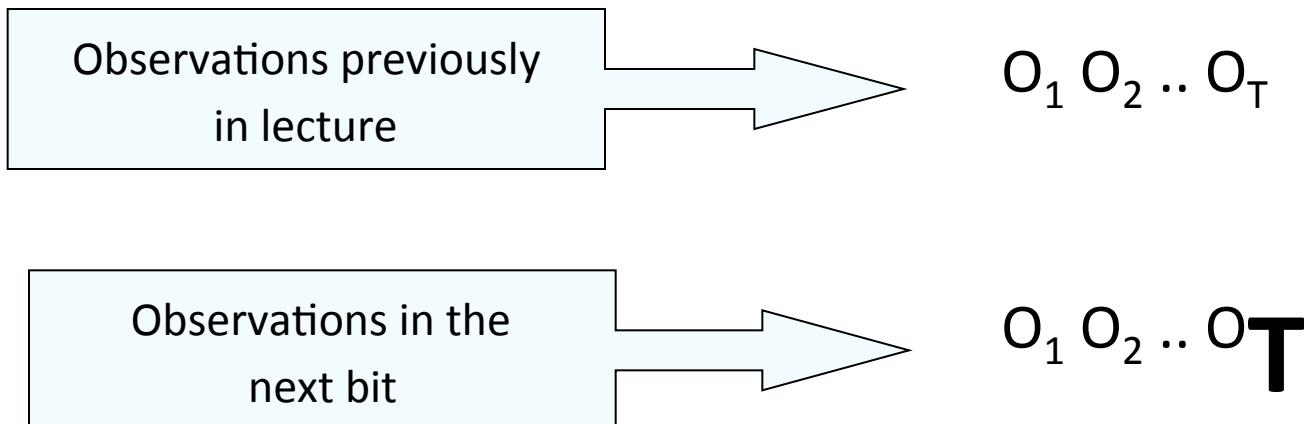
In practice: many levels of inference; not
one big jump.

HMMs are used and useful

But how do you design an HMM?

Occasionally, (e.g. in our robot example) it is reasonable to deduce the HMM from first principles.

But usually, especially in Speech or Genetics, it is better to infer it from large amounts of data. $O_1 O_2 \dots O_T$ with a big “T”.



Inferring an HMM

Remember, we have been doing things like

$$P(O_1 O_2 \dots O_T | \lambda)$$

That “ λ ” is the notation for our HMM parameters.

Now We have some observations and we want to estimate λ from them.

AS USUAL: We could use

(i) MAX LIKELIHOOD $\lambda = \underset{\lambda'}{\operatorname{argmax}} P(O_1 \dots O_T | \lambda')$

(ii) BAYES Work out $P(\lambda | O_1 \dots O_T)$

and then take $E[\lambda]$ or $\max_{\lambda} P(\lambda | O_1 \dots O_T)$

λ

Max likelihood HMM estimation

Define

$$\gamma_t(i) = P(q_t = S_i | O_1 O_2 \dots O_T, \lambda)$$

$$\varepsilon_t(i,j) = P(q_t = S_i \wedge q_{t+1} = S_j | O_1 O_2 \dots O_T, \lambda)$$

$\gamma_t(i)$ and $\varepsilon_t(i,j)$ can be computed efficiently $\forall i, j, t$

(Details in Rabiner paper)

$$\sum_{t=1}^{T-1} \gamma_t(i) = \text{Expected number of transitions out of state } i \text{ during the path}$$

$$\sum_{t=1}^{T-1} \varepsilon_t(i, j) = \text{Expected number of transitions from state } i \text{ to state } j \text{ during the path}$$

$$\gamma_t(i) = P(q_t = S_i | O_1 O_2 \dots O_T, \lambda)$$

$$\varepsilon_t(i, j) = P(q_t = S_i \wedge q_{t+1} = S_j | O_1 O_2 \dots O_T, \lambda)$$

$\sum_{t=1}^{T-1} \gamma_t(i)$ = expected number of transitions out of state i during path

$\sum_{t=1}^{T-1} \varepsilon_t(i, j)$ = expected number of transitions out of i and into j during path

HMM estimation

Notice $\frac{\sum_{t=1}^{T-1} \varepsilon_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)} = \begin{pmatrix} \text{expected frequency} \\ i \rightarrow j \\ \hline \text{expected frequency} \\ i \end{pmatrix}$
= Estimate of Prob(Next state $S_j | This state S_i$)

We can re-estimate

$$a_{ij} \leftarrow \frac{\sum \varepsilon_t(i, j)}{\sum \gamma_t(i)}$$

We can also re-estimate

$$b_j(O_k) \leftarrow \dots \quad (\text{See Rabiner})$$

We want a_{ij}^{new} = new estimate of $P(q_{t+1} = s_j \mid q_t = s_i)$

We want a_{ij}^{new} = new estimate of $P(q_{t+1} = s_j \mid q_t = s_i)$

$$= \frac{\text{Expected \# transitions } i \rightarrow j \mid \lambda^{old}, O_1, O_2, \dots, O_T}{\sum_{k=1}^N \text{Expected \# transitions } i \rightarrow k \mid \lambda^{old}, O_1, O_2, \dots, O_T}$$

We want $a_{ij}^{\text{new}} = \text{new estimate of } P(q_{t+1} = s_j \mid q_t = s_i)$

$$= \frac{\text{Expected \# transitions } i \rightarrow j \mid \lambda^{\text{old}}, O_1, O_2, \dots, O_T}{\sum_{k=1}^N \text{Expected \# transitions } i \rightarrow k \mid \lambda^{\text{old}}, O_1, O_2, \dots, O_T}$$

$$= \frac{\sum_{t=1}^T P(q_{t+1} = s_j, q_t = s_i \mid \lambda^{\text{old}}, O_1, O_2, \dots, O_T)}{\sum_{k=1}^N \sum_{t=1}^T P(q_{t+1} = s_k, q_t = s_i \mid \lambda^{\text{old}}, O_1, O_2, \dots, O_T)}$$

We want $a_{ij}^{\text{new}} = \text{new estimate of } P(q_{t+1} = s_j \mid q_t = s_i)$

$$\begin{aligned}
 &= \frac{\text{Expected \# transitions } i \rightarrow j \mid \lambda^{\text{old}}, O_1, O_2, \dots, O_T}{\sum_{k=1}^N \text{Expected \# transitions } i \rightarrow k \mid \lambda^{\text{old}}, O_1, O_2, \dots, O_T} \\
 &= \frac{\sum_{t=1}^T P(q_{t+1} = s_j, q_t = s_i \mid \lambda^{\text{old}}, O_1, O_2, \dots, O_T)}{\sum_{k=1}^N \sum_{t=1}^T P(q_{t+1} = s_k, q_t = s_i \mid \lambda^{\text{old}}, O_1, O_2, \dots, O_T)}
 \end{aligned}$$

$$\begin{aligned}
 &= \frac{S_{ij}}{\sum_{k=1}^N S_{ik}} \text{ where } S_{ij} = \sum_{t=1}^T P(q_{t+1} = s_j, q_t = s_i, O_1, \dots, O_T \mid \lambda^{\text{old}})
 \end{aligned}$$

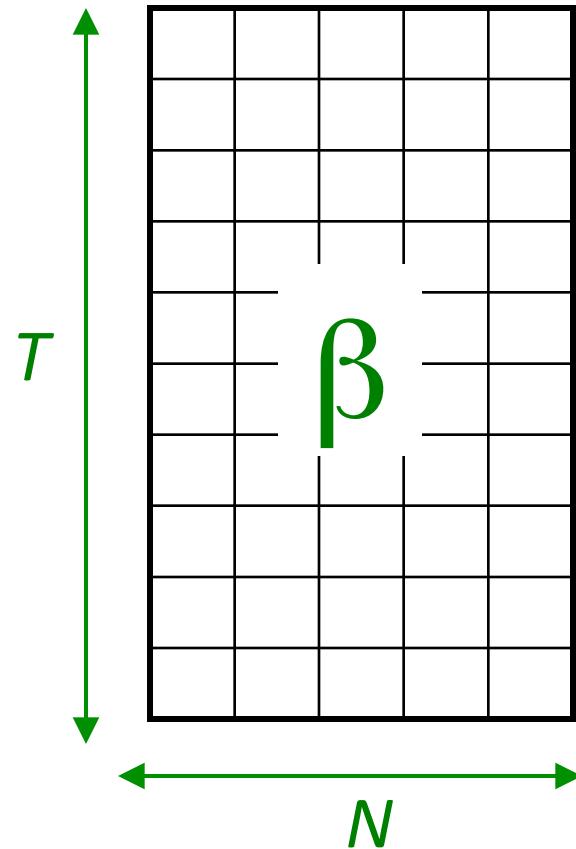
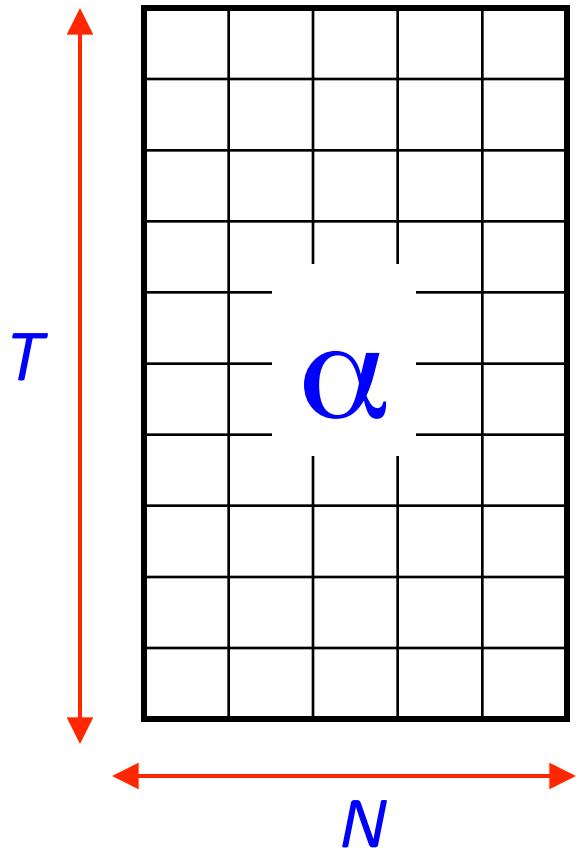
= What?

We want $a_{ij}^{\text{new}} = \text{new estimate of } P(q_{t+1} = s_j \mid q_t = s_i)$

$$\begin{aligned}
 &= \frac{\text{Expected \# transitions } i \rightarrow j \mid \lambda^{\text{old}}, O_1, O_2, \dots, O_T}{\sum_{k=1}^N \text{Expected \# transitions } i \rightarrow k \mid \lambda^{\text{old}}, O_1, O_2, \dots, O_T} \\
 &= \frac{\sum_{t=1}^T P(q_{t+1} = s_j, q_t = s_i \mid \lambda^{\text{old}}, O_1, O_2, \dots, O_T)}{\sum_{k=1}^N \sum_{t=1}^T P(q_{t+1} = s_k, q_t = s_i \mid \lambda^{\text{old}}, O_1, O_2, \dots, O_T)}
 \end{aligned}$$

$$\begin{aligned}
 &= \frac{S_{ij}}{\sum_{k=1}^N S_{ik}} \text{ where } S_{ij} = \sum_{t=1}^T P(q_{t+1} = s_j, q_t = s_i, O_1, \dots, O_T \mid \lambda^{\text{old}}) \\
 &\quad = a_{ij} \sum_{t=1}^T \alpha_t(i) \beta_{t+1}(j) b_j(O_{t+1})
 \end{aligned}$$

We want $a_{ij}^{\text{new}} = S_{ij} \Bigg/ \sum_{k=1}^N S_{ik}$ where $S_{ij} = a_{ij} \sum_{t=1}^T \alpha_t(i) \beta_{t+1}(j) b_j(O_{t+1})$



EM for HMMs

If we knew λ we could estimate EXPECTATIONS of quantities such as

- Expected number of times in state i
- Expected number of transitions $i \rightarrow j$

If we knew the quantities such as

- Expected number of times in state i
- Expected number of transitions $i \rightarrow j$

We could compute the MAX LIKELIHOOD estimate of

$$\lambda = \langle \{a_{ij}\}, \{b_i(j)\}, \pi_i \rangle$$

Roll on the EM Algorithm...

EM for HMMs

1. Get your observations $O_1 \dots O_T$
 2. Guess your first λ estimate $\lambda(0)$, $k=0$
 3. $k = k+1$
 4. Given $O_1 \dots O_T$, $\lambda(k)$ compute
$$\gamma_t(i), \varepsilon_t(i,j) \quad \forall 1 \leq t \leq T, \quad \forall 1 \leq i \leq N, \quad \forall 1 \leq j \leq N$$
 5. Compute expected freq. of state i , and expected freq. $i \rightarrow j$
 6. Compute new estimates of a_{ij} , $b_j(k)$, π_i accordingly. Call them $\lambda(k+1)$
 7. Goto 3, unless converged.
- **Also known (for the HMM case) as the BAUM-WELCH algorithm.**



Baum-Welch



- **Leonard Esau Baum** (August 23, 1931 – August 14, 2017) was an American mathematician. He graduated from [Harvard University](#) in 1953, and earned a Ph.D. in mathematics from Harvard in 1958, with a dissertation entitled *Derivations in Commutative Semi-Simple Banach Algebras*. At the time he wrote his work with Welch, he was working for the [Institute for Defense Analyses](#) in [Princeton, New Jersey](#).
- **Lloyd Richard Welch** (born September 28, 1927) is an American [information theorist](#) and applied mathematician, and co-inventor of the [Baum–Welch algorithm](#) and the [Berlekamp–Welch algorithm](#), also known as the Welch–Berlekamp algorithm. Welch received his B.S. in mathematics from the [University of Illinois](#), 1951, and Ph.D. in mathematics from the [California Institute of Technology](#), 1958.

Bad News

- There are lots of local minima

Good News

- The local minima are usually adequate models of the data.

Notice

- EM does not estimate the number of states. That must be given.
- Often, HMMs are forced to have some links with zero probability. This is done by setting $a_{ij}=0$ in initial estimate $\lambda(0)$
- Easy extension of everything seen today: HMMs with real valued outputs

Basic Operations in HMMs

For an observation sequence $O = O_1 \dots O_T$, the three basic HMM operations are:

Problem	Algorithm	Complexity ₊
<i>Evaluation:</i> Calculating $P(q_t=S_i O_1 O_2 \dots O_t)$	Forward-Backward	$O(TN^2)$
<i>Inference:</i> Computing $Q^* = \text{argmax}_Q P(Q O)$	Viterbi Decoding	$O(TN^2)$
<i>Learning:</i> Computing $\lambda^* = \text{argmax}_\lambda P(O \lambda)$	Baum-Welch (EM)	$O(TN^2)$

$T = \# \text{ timesteps}$, $N = \# \text{ states}$



What You Should Know

- What is an HMM ?
- Computing (and defining) $\alpha_t(i)$
- The Viterbi algorithm
- Outline of the EM algorithm
- Good: fairly thorough reading of Rabiner* from page 257 to page 266* [Up to but not including “IV. Types of HMMs”].

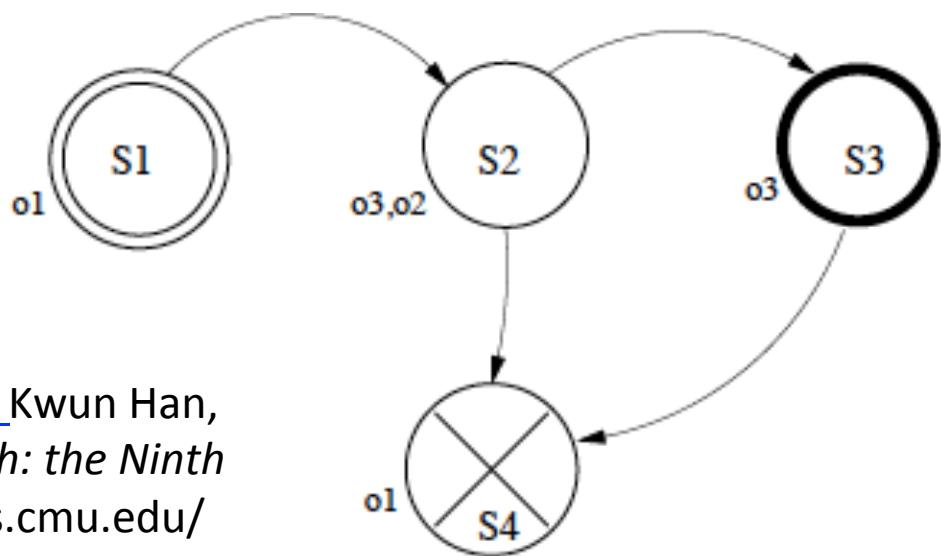
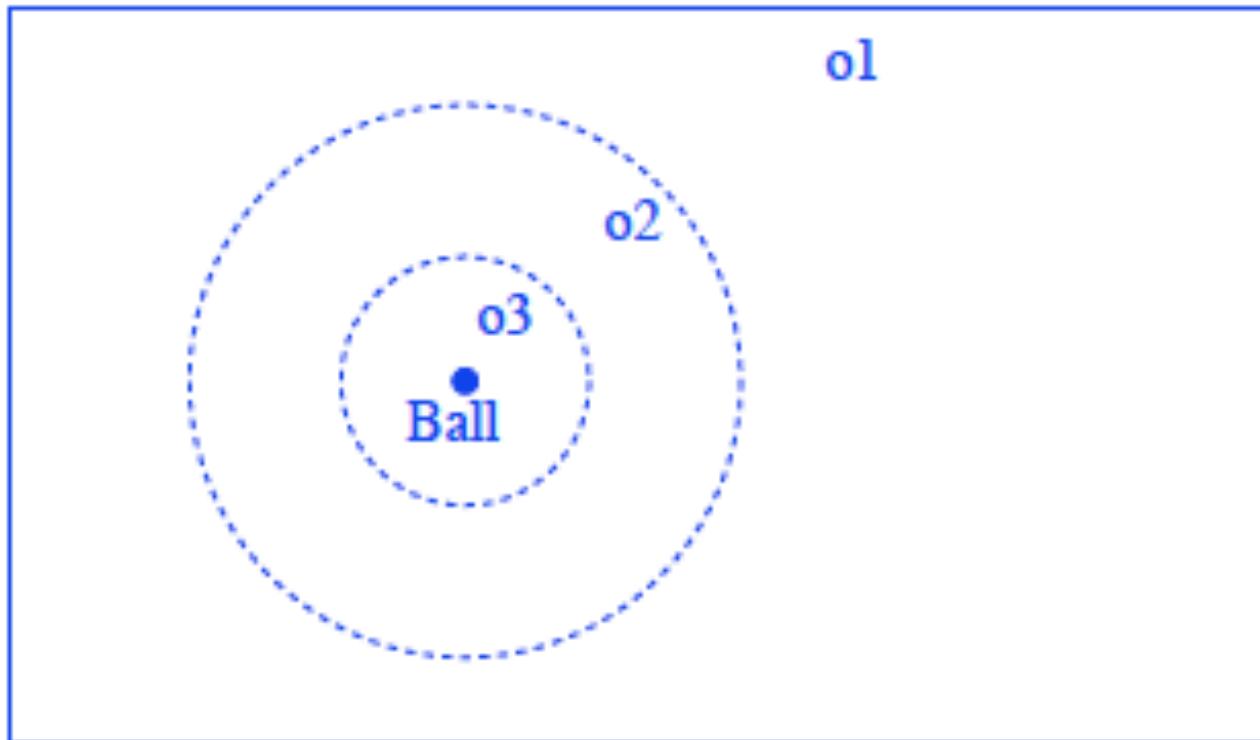
*L. R. Rabiner, "A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition," Proc. of the IEEE, Vol.77, No.2, pp. 257--286, 1989.

<http://ieeexplore.ieee.org/iel5/5/698/00018626.pdf?arnumber=18626>

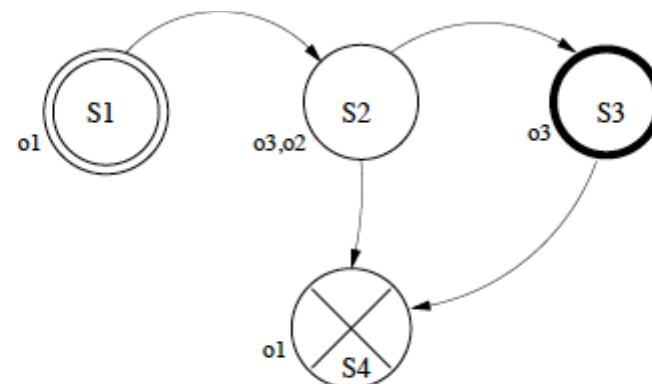
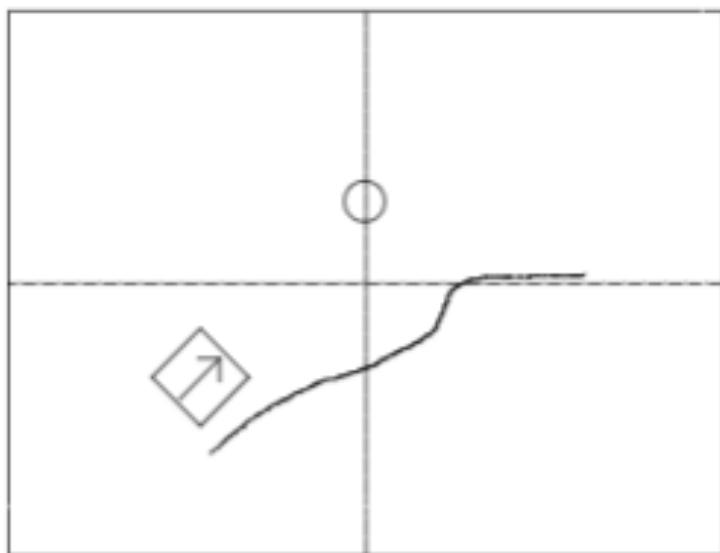
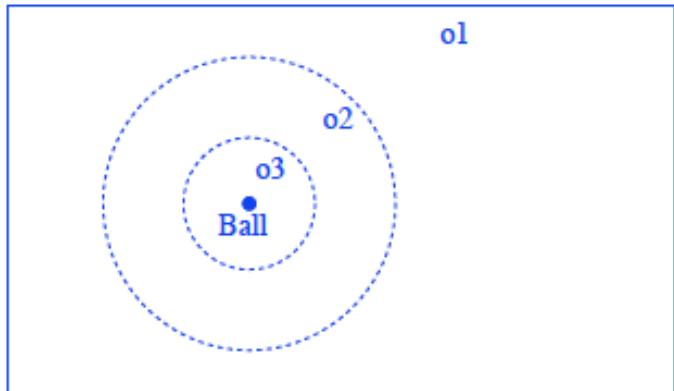
Applications of HMMs

Speech Recognition

- Observations: features from speech segment
- States: word being spoken
- State Transition Model: Language Model on word sequences
- Observation Model: Acoustic Model on speech sounds for specific words



[Automated Robot Behavior Recognition](#), Kwun Han,
and Manuela Veloso In *Robotics Research: the Ninth
International Symposium*, [http://www.cs.cmu.edu/
~mmv/papers/ijcai99-kwun.pdf](http://www.cs.cmu.edu/~mmv/papers/ijcai99-kwun.pdf)

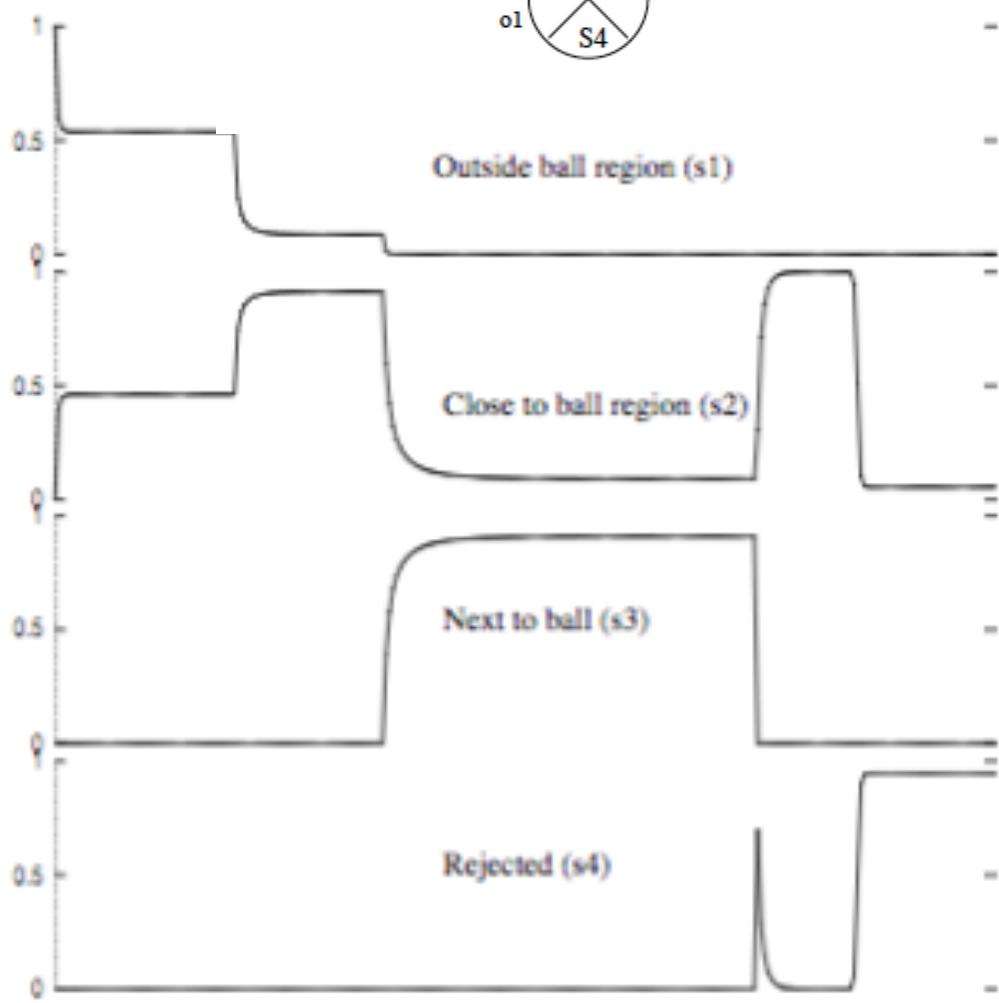


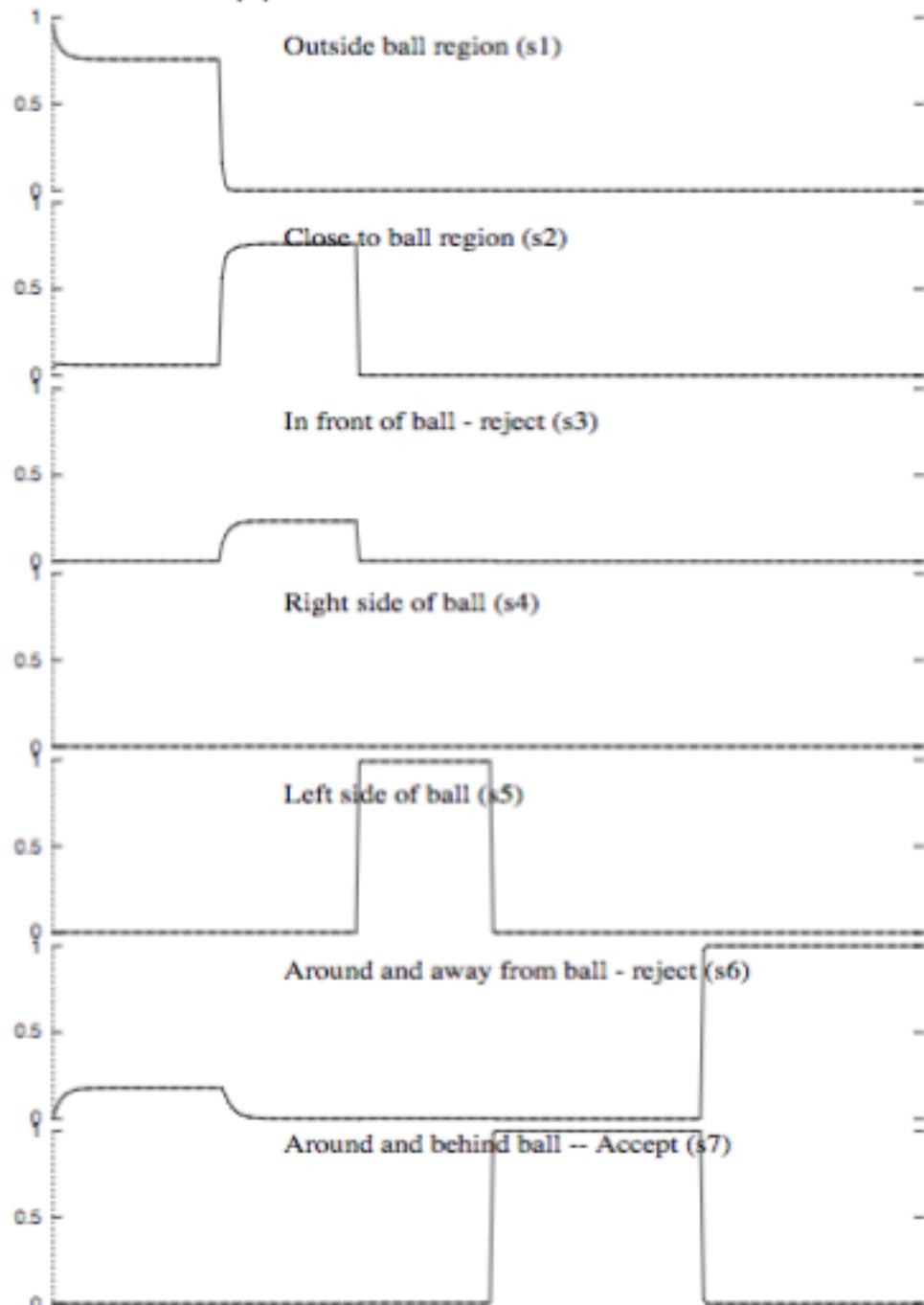
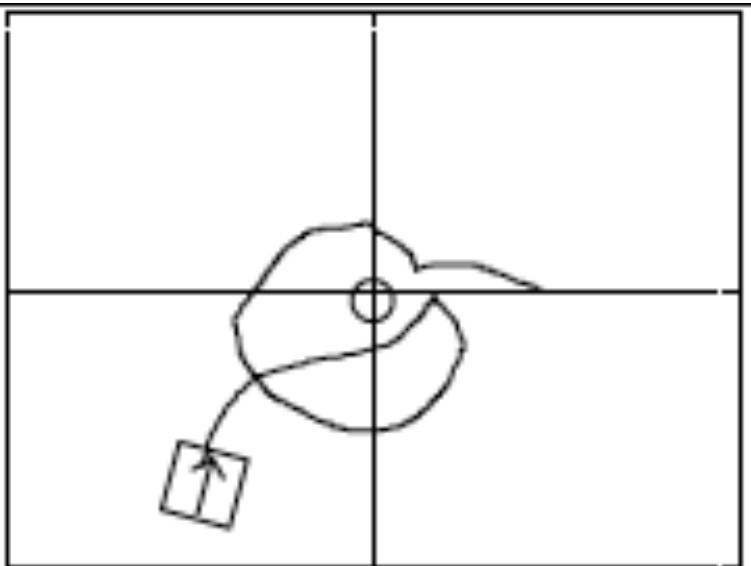
Outside ball region (s_1)

Close to ball region (s_2)

Next to ball (s_3)

Rejected (s_4)





(a) Go-Behind-Ball Behavior

Activity Recognition from Videos

- Observations: image features: position displacement
- States: Type of activity: stand, standing, walk left, walk right, sit, sitting



[FOCUS: A Generalized Method for Object Discovery for Robots that Observe and Interact with Humans,](#) Manuela Veloso, Paul Rybski, and Felix von Hundelshausen
In HRI 2006, <http://www.cs.cmu.edu/~mmv/papers/06hri-focus.pdf>

Part of Speech Tagging

- Observations: Words
- States: Part of Speech (noun, verb, adjective, etc.)
- State Transition Model: Grammar Model on POS sequences
- Observation Model: Language Model on Word | POS