

HMM MDP RL

Mike Stilman (and a very frustrated robot)
robot@cmu.edu

The Robotics Institute
Carnegie Mellon University

Machine Learning 10-701 Course Review, Fall 2005

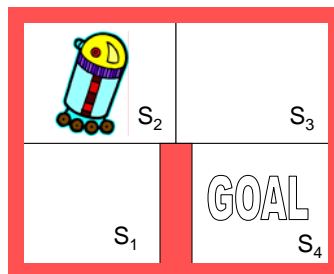
Copyright © 2005, Mike Stilman

Dec 14th, 2005

THE WORLD

A robot is trying to get to the goal

- Robot = R
- Goal = S_4
- 4 States



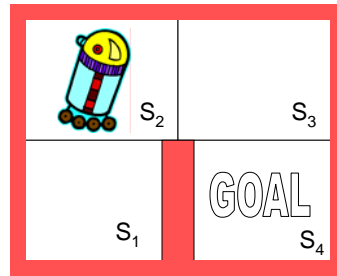
Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 2

Markov?

Can we represent this as a Markov System?

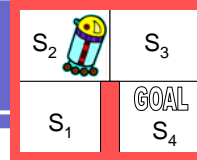
- Why?



Copyright © 2005, Mike Stilman

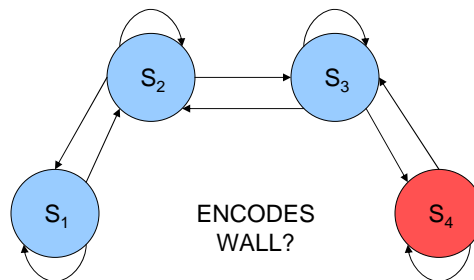
HMM MDP RL: Slide 3

Ok - Markov



Markov Model:

- Legitimate Transitions



Copyright © 2005, Mike Stilman

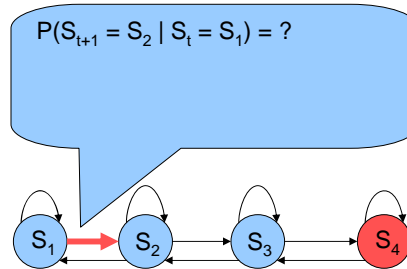
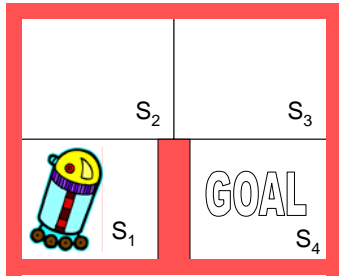
HMM MDP RL: Slide 4

Elementary Markov Learning

● Lets Watch the Robot

It Does Sequence:

$S_1 S_2 S_1 S_1 S_2 S_2 S_1$
 $S_2 S_3 S_2 S_2 S_2 S_2 S_3 S_4$
 $S_4 S_3 S_3 S_3 S_3$



Copyright © 2005, Mike Stilman

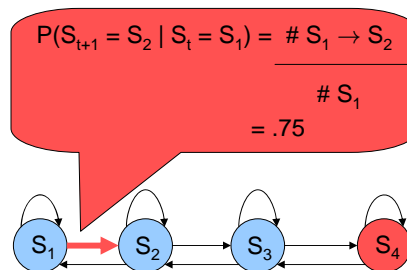
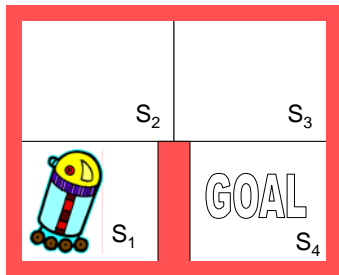
HMM MDP RL: Slide 5

Elementary Markov Learning

● Lets Watch the Robot

It Does Sequence:

$S_1 S_2 S_1 S_1 S_2 S_2 S_1$
 $S_2 S_3 S_2 S_2 S_2 S_2 S_3 S_4$
 $S_4 S_3 S_3 S_3 S_3$



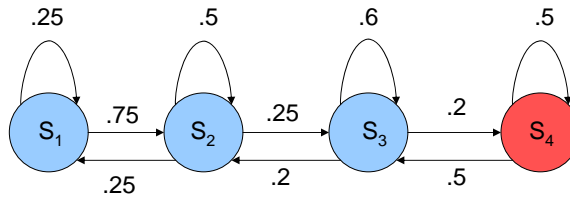
Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 6

Elementary Markov Learning

● Lets Watch the Robot

It Does Sequence:
 $S_1 S_2 S_1 S_1 S_2 S_2 S_1$
 $S_2 S_3 S_2 S_2 S_2 S_2 S_3 S_4$
 $S_4 S_3 S_3 S_3 S_3$



Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 7

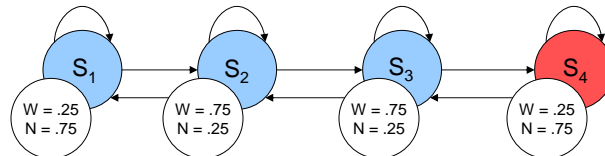
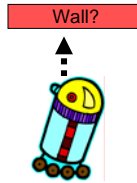
HMM: Observations



OK, we learned transition probabilities

Lets expand our model:

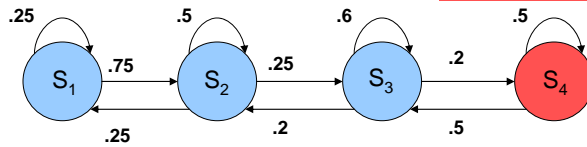
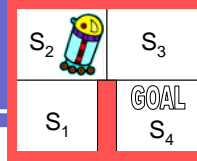
- Robots have sensors!
- Our robot can only look up
- $\frac{1}{4}$ of the time it's wrong
(hence frustrated)



Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 8

HMM: Transitions

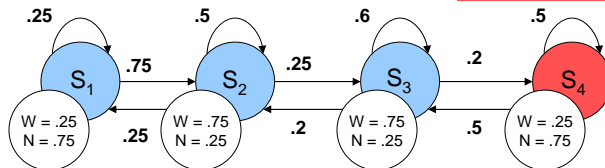
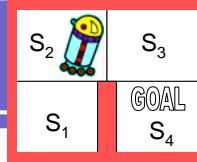


A

	s_1	s_2	s_3	s_4
s_1	.25	.75	0	0
s_2	.25	.5	.25	0
s_3	0	.2	.6	.2
s_4	0	0	.5	.5

Transition Matrix

HMM: Observations



A

	s_1	s_2	s_3	s_4
s_1	.25	.75	0	0
s_2	.25	.5	.25	0
s_3	0	.2	.6	.2
s_4	0	0	.5	.5

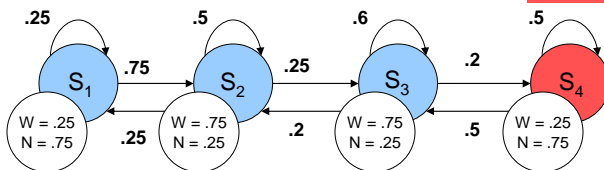
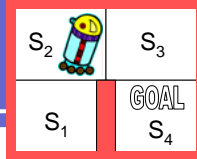
Transition Matrix

B

	W	N
s_1	.25	.75
s_2	.75	.25
s_3	.75	.25
s_4	.25	.75

Observation Matrix

HMM: λ Complete Model



A

	s_1	s_2	s_3	s_4
s_1	.25	.75	0	0
s_2	.25	.5	.25	0
s_3	0	.2	.6	.2
s_4	0	0	.5	.5

B

	W	N
s_1	.25	.75
s_2	.75	.25
s_3	.75	.25
s_4	.25	.75

π

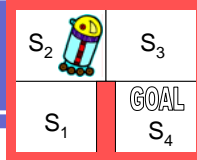
	π
s_1	.5
s_2	.3
s_3	.1
s_4	.1

Transition Matrix

Observation Matrix

Initial

HMM: Where am I?



- The robot has observed:
 - No Wall
 - Wall
 - Wall
 - No Wall
- What is the probability that the robot is now at the goal?

$P(q_4 = S_4 \mid N_1 W_2 W_3 N_4 \lambda) = ?$



HMM: Where am I?

FORWARD ALGORITHM

$$\alpha_t(i) = P(O_1 \dots O_t \mid q_t = S_i \mid \lambda)$$

- Initialization: $\alpha_1(i) = P(O_1 \mid q_1 = S_i \mid \lambda) = \pi(S_i) P(O_1 \mid S_i)$

	N	W	W	N
S ₁	.5 × .75 = .375			
S ₂	.3 × .25 = .075			
S ₃	.1 × .25 = .025			
S ₄	.1 × .75 = .075			

	S ₁	S ₂	S ₃	S ₄
S ₁	.25	.75	0	0
S ₂	.25	.5	.25	0
S ₃	0	.2	.6	.2
S ₄	0	0	.5	.5

	W	N
S ₁	.25	.75
S ₂	.75	.25
S ₃	.75	.25
S ₄	.25	.75

	π
S ₁	.5
S ₂	.3
S ₃	.1
S ₄	.1

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 13

HMM: Where am I?

FORWARD ALGORITHM

$$\alpha_t(i) = P(O_1 \dots O_t \mid q_t = S_i \mid \lambda)$$

- Induction: $\alpha_{t+1}(i) = [\sum_j \alpha_t(j) A_{ij}] B_j(O_{t+1})$

	N	W	W	N
S ₁	.375 → (.375 × .25 + .075 × .25) × .25			
S ₂	.075			
S ₃	.025			
S ₄	.075			

	S ₁	S ₂	S ₃	S ₄
S ₁	.25	.75	0	0
S ₂	.25	.5	.25	0
S ₃	0	.2	.6	.2
S ₄	0	0	.5	.5

	W	N
S ₁	.25	.75
S ₂	.75	.25
S ₃	.75	.25
S ₄	.25	.75

	π
S ₁	.5
S ₂	.3
S ₃	.1
S ₄	.1

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 14

HMM: Where am I?

FORWARD ALGORITHM

$$\alpha_t(i) = P(O_1 \dots O_t \mid q_t = S_i \mid \lambda)$$

- Induction: $\alpha_{t+1}(i) = [\sum_j \alpha_t(j) A_{ij}] B_j(O_{t+1})$

	N	W	W	N
S ₁	.375	.1125		
S ₂	.075	(.375 × .75 + .075 × .5 → .025 × .2) × .75		
S ₃	.025			
S ₄	.075			

	S ₁	S ₂	S ₃	S ₄
S ₁	.25	.75	0	0
S ₂	.25	.5	.25	0
S ₃	0	.2	.6	.2
S ₄	0	0	.5	.5

	W	N
S ₁	.25	.75
S ₂	.75	.25
S ₃	.75	.25
S ₄	.25	.75

	π
S ₁	.5
S ₂	.3
S ₃	.1
S ₄	.1

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 15

HMM: Where am I?

FORWARD ALGORITHM

$$\alpha_t(i) = P(O_1 \dots O_t \mid q_t = S_i \mid \lambda)$$

- Induction: $\alpha_{t+1}(i) = [\sum_j \alpha_t(j) A_{ij}] B_j(O_{t+1})$

	N	W	W	N
S ₁	.375	.1125		
S ₂	.075	.2428		
S ₃	.025	(.075 × .25 + .025 × .6 → .075 × .5) × .75		
S ₄	.075			

	S ₁	S ₂	S ₃	S ₄
S ₁	.25	.75	0	0
S ₂	.25	.5	.25	0
S ₃	0	.2	.6	.2
S ₄	0	0	.5	.5

	W	N
S ₁	.25	.75
S ₂	.75	.25
S ₃	.75	.25
S ₄	.25	.75

	π
S ₁	.5
S ₂	.3
S ₃	.1
S ₄	.1

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 16

HMM: Where am I?

FORWARD ALGORITHM

$$\alpha_t(i) = P(O_1 \dots O_t \mid q_t = S_i \mid \lambda)$$

- Induction: $\alpha_{t+1}(i) = [\sum_j \alpha_t(j) A_{ij}] B_j(O_{t+1})$

	N	W	W	N
S ₁	.375	.1125		
S ₂	.075	.2428		
S ₃	.025	.0534		
S ₄	.075	.0106		

	S ₁	S ₂	S ₃	S ₄
S ₁	.25	.75	0	0
S ₂	.25	.5	.25	0
S ₃	0	.2	.6	.2
S ₄	0	0	.5	.5

	W	N
S ₁	.25	.75
S ₂	.75	.25
S ₃	.75	.25
S ₄	.25	.75

	π
S ₁	.5
S ₂	.3
S ₃	.1
S ₄	.1

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 17

HMM: Where am I?

FORWARD ALGORITHM

$$\alpha_t(i) = P(O_1 \dots O_t \mid q_t = S_i \mid \lambda)$$

- Induction: $\alpha_{t+1}(i) = [\sum_j \alpha_t(j) A_{ij}] B_j(O_{t+1})$

	N	W	W	N
S ₁	.375	.1125	$(.1125 \times .25 + .2482 \times .25) \times .25$	
S ₂	.075	.2428		
S ₃	.025	.0534		
S ₄	.075	.0106		

	S ₁	S ₂	S ₃	S ₄
S ₁	.25	.75	0	0
S ₂	.25	.5	.25	0
S ₃	0	.2	.6	.2
S ₄	0	0	.5	.5

	W	N
S ₁	.25	.75
S ₂	.75	.25
S ₃	.75	.25
S ₄	.25	.75

	π
S ₁	.5
S ₂	.3
S ₃	.1
S ₄	.1

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 18

HMM: Where am I?

FORWARD ALGORITHM

$$\alpha_t(i) = P(O_1 \dots O_t \mid q_t = S_i \mid \lambda)$$

- Induction: $\alpha_{t+1}(i) = [\sum_j \alpha_t(j) A_{ij}] B_j(O_{t+1})$

	N	W	W	N
S ₁	.375	.1125	.0222	
S ₂	.075	.2428	.1623	
S ₃	.025	.0534	.0735	
S ₄	.075	.0106	.0040	

	S ₁	S ₂	S ₃	S ₄
S ₁	.25	.75	0	0
S ₂	.25	.5	.25	0
S ₃	0	.2	.6	.2
S ₄	0	0	.5	.5

	W	N
S ₁	.25	.75
S ₂	.75	.25
S ₃	.75	.25
S ₄	.25	.75

	π
S ₁	.5
S ₂	.3
S ₃	.1
S ₄	.1

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 19

HMM: Where am I?

FORWARD ALGORITHM

$$\alpha_t(i) = P(O_1 \dots O_t \mid q_t = S_i \mid \lambda)$$

- Induction: $\alpha_{t+1}(i) = [\sum_j \alpha_t(j) A_{ij}] B_j(O_{t+1})$

	N	W	W	N
S ₁	.375	.1125	.0222	.0346
S ₂	.075	.2428	.1623	.0281
S ₃	.025	.0534	.0735	.0217
S ₄	.075	.0106	.0040	.0125

	S ₁	S ₂	S ₃	S ₄
S ₁	.25	.75	0	0
S ₂	.25	.5	.25	0
S ₃	0	.2	.6	.2
S ₄	0	0	.5	.5

	W	N
S ₁	.25	.75
S ₂	.75	.25
S ₃	.75	.25
S ₄	.25	.75

	π
S ₁	.5
S ₂	.3
S ₃	.1
S ₄	.1

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 20

HMM: Where am I?

- We found $P(N_1 W_2 W_3 N_4 q_4 = S_4 | \lambda) = .0125$
- Is this the probability of being in S_4 given that we have seen N_1, W_2, W_3 and N_4 ?

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 21

HMM: Where am I?

- We found $P(N_1 W_2 W_3 N_4 q_4 = S_4 | \lambda) = .0125$
- Is this the probability of being in S_4 given that we have seen N_1, W_2, W_3 and N_4 ?

NO

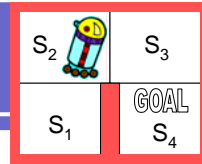
$$\begin{aligned} P(q_4 = S_4 | N_1 W_2 W_3 N_4 \lambda) &= \frac{P(N_1 W_2 W_3 N_4 q_4 = S_4 | \lambda)}{P(N_1 W_2 W_3 N_4)} \\ &= .0125 / (.0346 + .0281 + .0217 + .0125) \\ &= .129 \end{aligned}$$

13% in State 4
36% in State 1

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 22

HMM: How Did I Get There?



- The robot has observed:
 - No Wall
 - Wall
 - Wall
 - No Wall

What is the most likely sequence of states that occurred?



Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 23

HMM: How Did I Get There?

VITERBI ALGORITHM

$$\delta_t(i) = \max_{q_1 \dots q_t} P(q_1 \dots q_t = i, O_1 \dots O_t | \lambda)$$

- Initialization: $\alpha_1(i) = P(O_1 | q_1 = S_i) = \pi(S_i)$

	N	W	W	N
S ₁	.5 × .75 = .375			
S ₂	.3 × .25 = .075			
S ₃	.1 × .25 = .025			
S ₄	.1 × .75 = .075			

A

	S ₁	S ₂	S ₃	S ₄
S ₁	.25	.75	0	0
S ₂	.25	.5	.25	0
S ₃	0	.2	.6	.2
S ₄	0	0	.5	.5

B

	W	N
S ₁	.25	.75
S ₂	.75	.25
S ₃	.75	.25
S ₄	.25	.75

π

	π
S ₁	.5
S ₂	.3
S ₃	.1
S ₄	.1

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 24

HMM: How Did I Get There?

VITERBI ALGORITHM

$$\delta_t(i) = \max_{q_1 \dots q_t} P(q_1 \dots q_t = i, O_1 \dots O_t | \lambda)$$

- Induction: $\delta_t(i) = [\max_j \delta_t(i) A_{ij}] B_j(O_{t+1})$

	N	W	W	N
S ₁	.375	→ max(.375 × .25 .075 × .25) × .25		
S ₂	.075			
S ₃	.025			
S ₄	.075			

	S ₁	S ₂	S ₃	S ₄
S ₁	.25	.75	0	0
S ₂	.25	.5	.25	0
S ₃	0	.2	.6	.2
S ₄	0	0	.5	.5

	W	N
S ₁	.25	.75
S ₂	.75	.25
S ₃	.75	.25
S ₄	.25	.75

	π
S ₁	.5
S ₂	.3
S ₃	.1
S ₄	.1

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 25

HMM: How Did I Get There?

VITERBI ALGORITHM

$$\delta_t(i) = \max_{q_1 \dots q_t} P(q_1 \dots q_t = i, O_1 \dots O_t | \lambda)$$

- Induction: $\delta_t(i) = [\max_j \delta_t(i) A_{ij}] B_j(O_{t+1})$

	N	W	W	N
S ₁	.375	→ .024		
S ₂	.075	↘ max (.375 × .75 .075 × .5 .025 × .2) × .75		
S ₃	.025			
S ₄	.075			

	S ₁	S ₂	S ₃	S ₄
S ₁	.25	.75	0	0
S ₂	.25	.5	.25	0
S ₃	0	.2	.6	.2
S ₄	0	0	.5	.5

	W	N
S ₁	.25	.75
S ₂	.75	.25
S ₃	.75	.25
S ₄	.25	.75

	π
S ₁	.5
S ₂	.3
S ₃	.1
S ₄	.1

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 26

HMM: How Did I Get There?

VITERBI ALGORITHM

$$\delta_t(i) = \max_{q_1 \dots q_t} P(q_1 \dots q_t = i, O_1 \dots O_t | \lambda)$$

- Induction: $\delta_t(i) = [\max_j \delta_t(i) A_{ij}] B_j(O_{t+1})$

	N	W	W	N
S ₁	.375	.024		
S ₂	.075	.211		
S ₃	.025	max (.075 × .25 .025 × .6 .075 × .5) × .75		
S ₄	.075			

	S ₁	S ₂	S ₃	S ₄
S ₁	.25	.75	0	0
S ₂	.25	.5	.25	0
S ₃	0	.2	.6	.2
S ₄	0	0	.5	.5

	W	N
S ₁	.25	.75
S ₂	.75	.25
S ₃	.75	.25
S ₄	.25	.75

	π
S ₁	.5
S ₂	.3
S ₃	.1
S ₄	.1

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 27

HMM: How Did I Get There?

VITERBI ALGORITHM

$$\delta_t(i) = \max_{q_1 \dots q_t} P(q_1 \dots q_t = i, O_1 \dots O_t | \lambda)$$

- Induction: $\delta_t(i) = [\max_j \delta_t(i) A_{ij}] B_j(O_{t+1})$

	N	W	W	N
S ₁	.375	.024		
S ₂	.075	.211		
S ₃	.025	.028		
S ₄	.075	.009		

	S ₁	S ₂	S ₃	S ₄
S ₁	.25	.75	0	0
S ₂	.25	.5	.25	0
S ₃	0	.2	.6	.2
S ₄	0	0	.5	.5

	W	N
S ₁	.25	.75
S ₂	.75	.25
S ₃	.75	.25
S ₄	.25	.75

	π
S ₁	.5
S ₂	.3
S ₃	.1
S ₄	.1

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 28

HMM: How Did I Get There?

VITERBI ALGORITHM

$$\delta_t(i) = \max_{q_1 \dots q_t} P(q_1 \dots q_t = i, O_1 \dots O_t | \lambda)$$

- Induction: $\delta_t(i) = [\max_j \delta_t(j) A_{ij}] B_i(O_{t+1})$

	N	W	W	N
S ₁	.375	→ .024	↗ .013	
S ₂	.075	↘ .211	→ .079	
S ₃	.025	↗ .028	↘ .040	
S ₄	.075	↘ .009	↘ .001	

	S ₁	S ₂	S ₃	S ₄
S ₁	.25	.75	0	0
S ₂	.25	.5	.25	0
S ₃	0	.2	.6	.2
S ₄	0	0	.5	.5

	W	N
S ₁	.25	.75
S ₂	.75	.25
S ₃	.75	.25
S ₄	.25	.75

	π
S ₁	.5
S ₂	.3
S ₃	.1
S ₄	.1

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 29

HMM: How Did I Get There?

VITERBI ALGORITHM

$$\delta_t(i) = \max_{q_1 \dots q_t} P(q_1 \dots q_t = i, O_1 \dots O_t | \lambda)$$

- Induction: $\delta_t(i) = [\max_j \delta_t(j) A_{ij}] B_i(O_{t+1})$

	N	W	W	N
S ₁	.375	→ .024	↗ .013	↗ .015
S ₂	.075	↘ .211	→ .079	→ .005
S ₃	.025	↗ .028	↘ .040	→ .006
S ₄	.075	↘ .009	↘ .001	↘ .006

	S ₁	S ₂	S ₃	S ₄
S ₁	.25	.75	0	0
S ₂	.25	.5	.25	0
S ₃	0	.2	.6	.2
S ₄	0	0	.5	.5

	W	N
S ₁	.25	.75
S ₂	.75	.25
S ₃	.75	.25
S ₄	.25	.75

	π
S ₁	.5
S ₂	.3
S ₃	.1
S ₄	.1

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 30

HMM: How Did I Get There?

VITERBI ALGORITHM

$$\delta_t(i) = \max_{q_1 \dots q_t} P(q_1 \dots q_t = i, O_1 \dots O_t | \lambda)$$

- Now, Just Trace it Back!

	N	W	W	N
S ₁	.375	.024	.013	.015
S ₂	.075	.211	.079	.005
S ₃	.025	.028	.040	.006
S ₄	.075	.009	.001	.006

	S ₁	S ₂	S ₃	S ₄
S ₁	.25	.75	0	0
S ₂	.25	.5	.25	0
S ₃	0	.2	.6	.2
S ₄	0	0	.5	.5

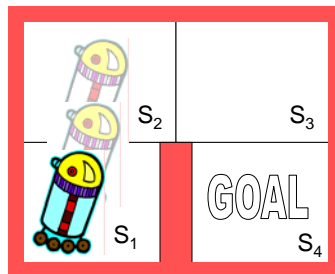
	W	N
S ₁	.25	.75
S ₂	.75	.25
S ₃	.75	.25
S ₄	.25	.75

	π
S ₁	.5
S ₂	.3
S ₃	.1
S ₄	.1

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 31

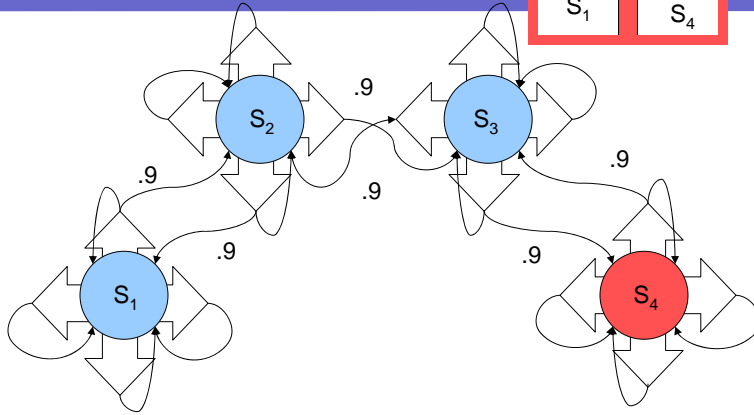
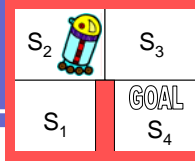
Robot's POLICY is...not great



Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 32

MDPs to the Rescue!

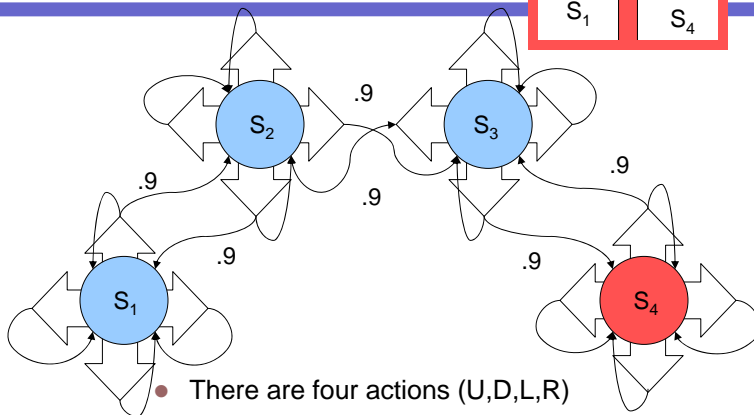
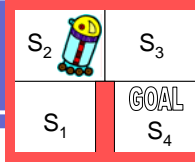


Assume the world is fully observable

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 33

MDPs to the Rescue!



- There are four actions (U,D,L,R)
- Moving into a wall makes the robot stay
- Moving to free space works 90% of the time. (Otherwise robot stays)

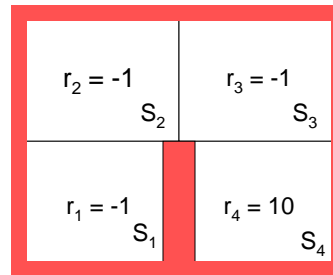
Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 34

Lets be Optimal

Find the optimal policy when:

- Rewards for leaving a state (r)
- Discount factor $\gamma = .9$
- **What is $J^*(S_i)$?**



Copyright © 2005, Mike Stilman

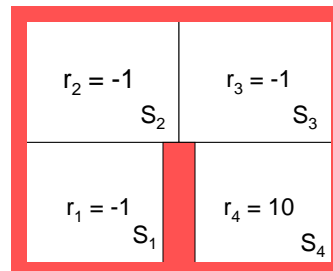
HMM MDP RL: Slide 35

Lets be Optimal

Find the optimal policy when:

- Rewards for leaving a state (r)
- Discount factor $\gamma = .9$
- **What is $J^*(S_i)$?**

$$J^*(S_i) = r_i + \gamma \max_a [\sum_j P(j | i, a) J^*(j)]$$



Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 36

Value Iteration $J^*(S_i) = r_i + \gamma \max_a [\sum_j P(j | i,a) J^*(j)]$

Iterative Local Optimization: Implements Dynamic Programming

$r_2 = -1$	$r_3 = -1$
$r_1 = -1$	$r_4 = 10$

	S_1	S_2	S_3	S_4
$J_0(S_i)$	0	0	0	0
$J_1(S_i)$				
$J_2(S_i)$				
$J_3(S_i)$				
$J_4(S_i)$				

Initialize to 0 – or something meaningful from policy

How do we update?

Value Iteration $J^*(S_i) = r_i + \gamma \max_a [\sum_j P(j | i,a) J^*(j)]$

Iterative Local Optimization: Implements Dynamic Programming

$r_2 = -1$	$r_3 = -1$
$r_1 = -1$	$r_4 = 10$

	S_1	S_2	S_3	S_4
$J_0(S_i)$	0	0	0	0
$J_1(S_i)$	-1			
$J_2(S_i)$				
$J_3(S_i)$				
$J_4(S_i)$				

$$J_{t+1}(S_i) = r_i + \gamma \max_a [\sum_j P(j | i,a) J_t(j)]$$

$$J_1(S_1) = -1 + .9 \max_a [1 \times 0, 1 \times 0, 1 \times 0, .9 \times 0 + .1 \times 0]$$

Value Iteration $J^*(S_i) = r_i + \gamma \max_a [\sum_j P(j | i, a) J^*(j)]$

Iterative Local Optimization: Implements Dynamic Programming

$r_2 = -1$	$r_3 = -1$
$r_1 = -1$	$r_4 = 10$

	S_1	S_2	S_3	S_4
$J_0(S_i)$	0	0	0	0
$J_1(S_i)$	-1	-1	-1	10
$J_2(S_i)$				
$J_3(S_i)$				
$J_4(S_i)$				

$$J_{t+1}(S_i) = r_i + \gamma \max_a [\sum_j P(j | i, a) J_t(j)]$$

$$J_1(S_1) = -1 + .9 \max_a [1 \times 0, 1 \times 0, 1 \times 0, .9 \times 0 + .1 \times 0]$$

...

$$J_1(S_4) = 10 + .9 \max_a [1 \times 0, 1 \times 0, 1 \times 0, .9 \times 0 + .1 \times 0]$$

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 39

Value Iteration $J^*(S_i) = r_i + \gamma \max_a [\sum_j P(j | i, a) J^*(j)]$

Iterative Local Optimization: Implements Dynamic Programming

$r_2 = -1$	$r_3 = -1$
$r_1 = -1$	$r_4 = 10$

	S_1	S_2	S_3	S_4
$J_0(S_i)$	0	0	0	0
$J_1(S_i)$	-1	-1	-1	10
$J_2(S_i)$	-1.9	-1.9	7.0	19
$J_3(S_i)$				
$J_4(S_i)$				

$$J_{t+1}(S_i) = r_i + \gamma \max_a [\sum_j P(j | i, a) J_t(j)]$$

$$J_2(S_1) = -1 + .9 \max_a [1 \times -1, 1 \times -1, 1 \times -1, .9 \times -1 + .1 \times -1]$$

...

$$J_2(S_3) = -1 + .9 \max_a [1 \times -1, 1 \times -1, 1 \times -1, .9 \times 10 + .1 \times -1]$$

$$J_2(S_4) = 10 + .9 \max_a [1 \times 10, 1 \times 10, 1 \times 10, .9 \times -1 + .1 \times 10]$$

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 40

Value Iteration $J^*(S_i) = r_i + \gamma \max_a [\sum_j P(j | i, a) J^*(j)]$

Iterative Local Optimization: Implements Dynamic Programming

$r_2 = -1$	$r_3 = -1$
$r_1 = -1$	$r_4 = 10$

	S_1	S_2	S_3	S_4
$J_0(S_i)$	0	0	0	0
$J_1(S_i)$	-1	-1	-1	10
$J_2(S_i)$	-1.9	-1.9	7.0	19
$J_3(S_i)$	-2.7	4.5	15.0	27.1
$J_4(S_i)$				

$$J_{t+1}(S_i) = r_i + \gamma \max_a [\sum_j P(j | i, a) J_t(j)]$$

Repeat until convergence...

Value Iteration $J^*(S_i) = r_i + \gamma \max_a [\sum_j P(j | i, a) J^*(j)]$

Iterative Local Optimization: Implements Dynamic Programming

$r_2 = -1$	$r_3 = -1$
$r_1 = -1$	$r_4 = 10$

	S_1	S_2	S_3	S_4
$J_0(S_i)$	0	0	0	0
$J_1(S_i)$	-1	-1	-1	10
$J_2(S_i)$	-1.9	-1.9	7.0	19
$J_3(S_i)$	-2.7	4.5	15.0	27.1
$J_4(S_i)$	2.4	11.6	22.3	34.39

$$J_{t+1}(S_i) = r_i + \gamma \max_a [\sum_j P(j | i, a) J_t(j)]$$

Repeat until convergence...

Value Iteration $J^*(S_i) = r_i + \gamma \max_a [\sum_j P(j | i, a) J^*(j)]$

Iterative Local Optimization: Implements Dynamic Programming

$r_2 = -1$	$r_3 = -1$
$r_1 = -1$	$r_4 = 10$

	S_1	S_2	S_3	S_4
$J_0(S_i)$	0	0	0	0
$J_1(S_i)$	-1	-1	-1	10
$J_2(S_i)$	-1.9	-1.9	7.0	19
$J_3(S_i)$	-2.7	4.5	15.0	27.1
$J_4(S_i)$	2.4	11.6	22.3	34.39

$$J_{t+1}(S_i) = r_i + \gamma \max_a [\sum_j P(j | i, a) J_t(j)]$$

To get policy, simply do argmax_a !

$$\pi_{t+1}(S_i) = \operatorname{argmax}_a [\sum_j P(j | i, a) J_t(j)]$$

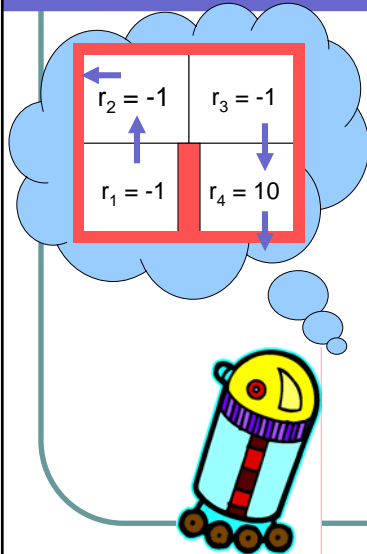
Downside...



	S_1	S_2	S_3	S_4
$J_0(S_i)$	0	0	0	0
$J_1(S_i)$	-1	-1	-1	10
$J_2(S_i)$	-1.9	-1.9	7.0	19
$J_3(S_i)$	-2.7	4.5	15.0	27.1
$J_4(S_i)$	2.4	11.6	22.3	34.39

Too Slow! I'm smarter than that.

Policy Iteration



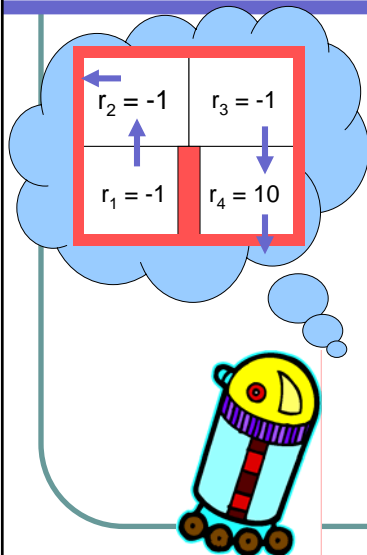
← Initial Policy

- 1) Find Value Function $J^*(S_i)$
 - Value Iteration (No Max)
 - Matrix Inversion

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 45

Policy Iteration



← Initial Policy

- 1) Find Value Function $J^*(S_i)$

$$J^*(S_4) = 10 \sum_{i=0}^{\infty} \gamma^i = 10 \times 1/(1-\gamma) = 100$$

$$J^*(S_3) = -1 + .9 [.9 J^*(S_4) + .1 J^*(S_3)]$$

$$J^*(S_3) = -1 + 81 + .09 J^*(S_3)$$

$$J^*(S_3) = 87.9$$

Analogously...

$$J^*(S_2) = -1 \sum_{i=0}^{\infty} \gamma^i = -10$$

$$J^*(S_1) = -10$$

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 46

Policy Iteration

$r_2 = -1$	$r_3 = -1$
$r_1 = -1$	$r_4 = 10$



- 1) Find Value Function $J^*(S_i)$
 - $J^*(S_4) = 100$
 - $J^*(S_3) = 87.9$
 - $J^*(S_2) = -10$
 - $J^*(S_1) = -10$
- 2) Compute *Improved* Policy

$$\pi_1(S_i) = \operatorname{argmax}_a [r_i + \gamma \sum_j P^a_{ij} J^*(S_j)]$$

$$\pi_1(S_2) = \operatorname{argmax}_a \left[\begin{array}{l} \leftarrow -1 + .9 (1 \times -10), \\ \rightarrow -1 + .9 (.9 \times 87.9 + .1 \times -10) \end{array} \right] = \operatorname{argmax}_a [-10, 69.3] = \text{RIGHT!} \rightarrow$$

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 47

Policy Iteration

$r_2 = -1$	$r_3 = -1$
$r_1 = -1$	$r_4 = 10$



← New Policy

- 1) Find Value Function $J^*(S_i)$
 - $J^*(S_4) = 100$
 - $J^*(S_3) = 87.9$
 - $J^*(S_2) = -10$
 - $J^*(S_1) = -10$
- 2) Compute *Improved* Policy

$$\pi_1(S_i) = \operatorname{argmax}_a [r_i + \gamma \sum_j P^a_{ij} J^*(S_j)]$$

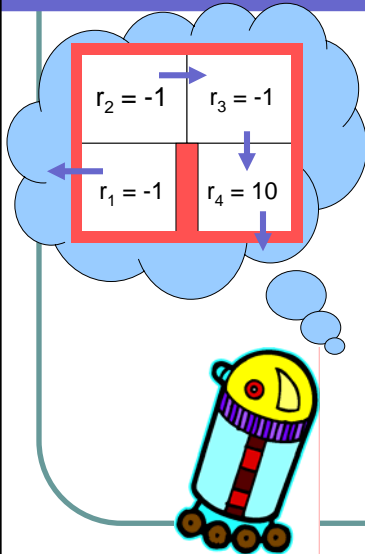
$$\begin{aligned} \pi_1(S_3) &= \text{DOWN} \\ \pi_1(S_4) &= \text{DOWN/LEFT/RIGHT} \\ \pi_1(S_1) &= \operatorname{argmax}_a [-10, -10] !!! \end{aligned}$$

What do we do?

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 48

Policy Iteration



← New Policy

- 1) Find Value Function $J^*(S_i)$
 - $J^*(S_4) = 100$
 - $J^*(S_3) = 87.9$
 - $J^*(S_2) = -10$
 - $J^*(S_1) = -10$
- 2) Compute *Improved Policy*

$$\pi_1(S_i) = \operatorname{argmax}_a [r_i + \gamma \sum_j P^a_{ij} J^*(S_j)]$$

$$\pi_1(S_1) = \operatorname{argmax}_a [-10, -10] !$$

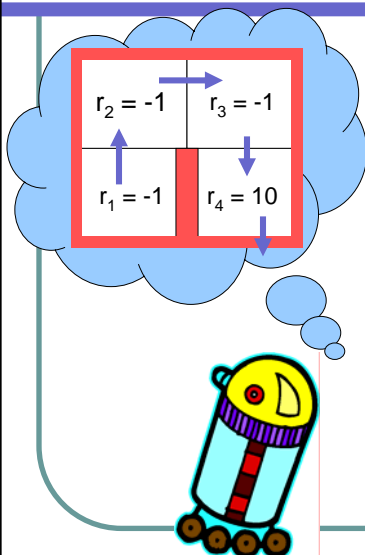
What do we do? Randomly choose left?

Then still bad policy, go back to step 1

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 49

Policy Iteration



← New Policy

- 1) Find Value Function $J^*(S_i)$
 - $J^*(S_4) = 100$
 - $J^*(S_3) = 87.9$
 - $J^*(S_2) = -10$
 - $J^*(S_1) = -10$
- 2) Compute *Improved Policy*

$$\pi_1(S_i) = \operatorname{argmax}_a [r_i + \gamma \sum_j P^a_{ij} J^*(S_j)]$$

$$\pi_1(S_1) = \operatorname{argmax}_a [-10, -10] !$$

What do we do? Randomly choose up?

Both cases must converge to this solution!

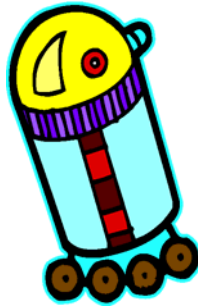
Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 50

We assumed a model for $P(j|i,a)$

What do we do if such a model does not exist?

- Make one
- Q-Learning

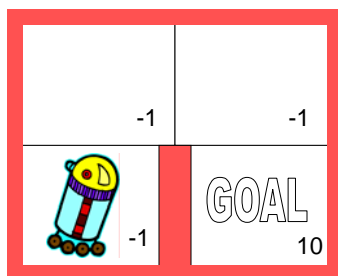


Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 51

Q-Learning

$$Q^{t+1}(S_i,a) \leftarrow \alpha [r_i + \gamma \max_{a^1} Q^t(S_i,a^1)] + (1-\alpha) Q^t(S_i,a)$$



$\alpha = .7$

Q-Table

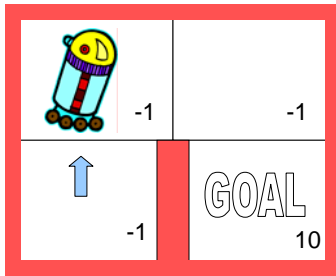
	↑	↓	←	→
S ₁	0	0	0	0
S ₂	0	0	0	0
S ₃	0	0	0	0
S ₄	0	0	0	0

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 52

Q-Learning

$$Q^{t+1}(S_i, a) \leftarrow \alpha [r_i + \gamma \max_{a^1} Q^t (S_j, a^1)] + (1-\alpha) Q^t(S_i, a)$$



$$Q^{\text{est}}(S_1, \uparrow) = .7(-1 + .9 \max(0, 0, 0, 0)) + .3 \times 0$$

Q-Table

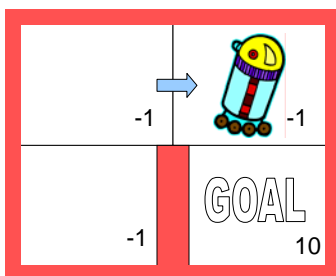
	↑	↓	←	→
S ₁	-.7	0	0	0
S ₂	0	0	0	0
S ₃	0	0	0	0
S ₄	0	0	0	0

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 53

Q-Learning

$$Q^{t+1}(S_i, a) \leftarrow \alpha [r_i + \gamma \max_{a^1} Q^t (S_j, a^1)] + (1-\alpha) Q^t(S_i, a)$$



$$Q^{\text{est}}(S_2, \rightarrow) = .7(-1 + .9 \max(0, 0, 0, 0)) + .3 \times 0$$

Q-Table

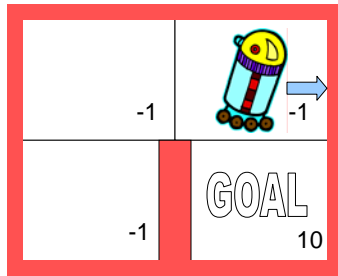
	↑	↓	←	→
S ₁	-.7	0	0	0
S ₂	0	0	0	-.7
S ₃	0	0	0	0
S ₄	0	0	0	0

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 54

Q-Learning

$$Q^{t+1}(S_i, a) \leftarrow \alpha [r_i + \gamma \max_{a^1} Q^t(S_j, a^1)] + (1-\alpha) Q^t(S_i, a)$$



$$Q^{\text{est}}(S_3, \rightarrow) = .7(-1 + .9 \max(0, 0, 0, 0)) + .3 \times 0$$

Q-Table

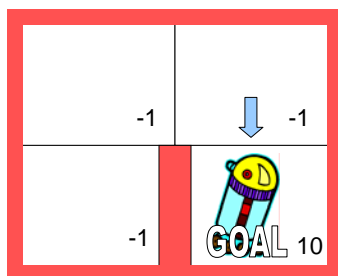
	↑	↓	←	→
S ₁	-0.7	0	0	0
S ₂	0	0	0	-0.7
S ₃	0	0	0	-0.7
S ₄	0	0	0	0

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 55

Q-Learning

$$Q^{t+1}(S_i, a) \leftarrow \alpha [r_i + \gamma \max_{a^1} Q^t(S_j, a^1)] + (1-\alpha) Q^t(S_i, a)$$



$$Q^{\text{est}}(S_3, \downarrow) = .7(-1 + .9 \max(0, 0, 0, 0)) + .3 \times 0$$

Q-Table

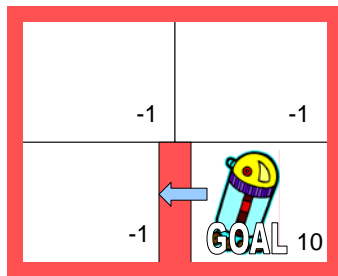
	↑	↓	←	→
S ₁	-0.7	0	0	0
S ₂	0	0	0	-0.7
S ₃	0	-0.7	0	-0.7
S ₄	0	0	0	0

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 56

Q-Learning

$$Q^{t+1}(S_i, a) \leftarrow \alpha [r_i + \gamma \max_{a^1} Q^t(S_j, a^1)] + (1-\alpha) Q^t(S_i, a)$$



$$Q^{\text{est}}(S_4, \leftarrow) = .7(10 + .9 \max(0, 0, 0, 0)) + .3 \times 0$$

Q-Table

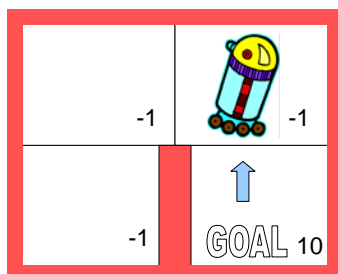
	↑	↓	←	→
S ₁	-0.7	0	0	0
S ₂	0	0	0	-0.7
S ₃	0	-0.7	0	-0.7
S ₄	0	0	7	0

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 57

Q-Learning

$$Q^{t+1}(S_i, a) \leftarrow \alpha [r_i + \gamma \max_{a^1} Q^t(S_j, a^1)] + (1-\alpha) Q^t(S_i, a)$$



$$Q^{\text{est}}(S_4, \uparrow) = .7(10 + .9 \max(0, -0.7, 0, -0.7)) + .3 \times 0$$

Q-Table

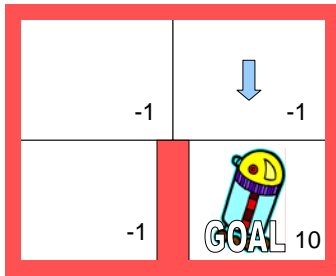
	↑	↓	←	→
S ₁	-0.7	0	0	0
S ₂	0	0	0	-0.7
S ₃	0	-0.7	0	-0.7
S ₄	7	0	7	0

Copyright © 2005, Mike Stilman

HMM MDP RL: Slide 58

Q-Learning

$$Q^{t+1}(S_i, a) \leftarrow \alpha [r_i + \gamma \max_{a^1} Q^t(S_j, a^1)] + (1-\alpha) Q^t(S_i, a)$$



$$Q^{\text{est}}(S_3, \downarrow) = .7(-1 + .9 \max(7, 0, 7, 0)) + .3 \times -7$$

Q-Table

	↑	↓	←	→
S ₁	-7	0	0	0
S ₂	0	0	0	-7
S ₃	0	3.5	0	-7
S ₄	7	0	7	0