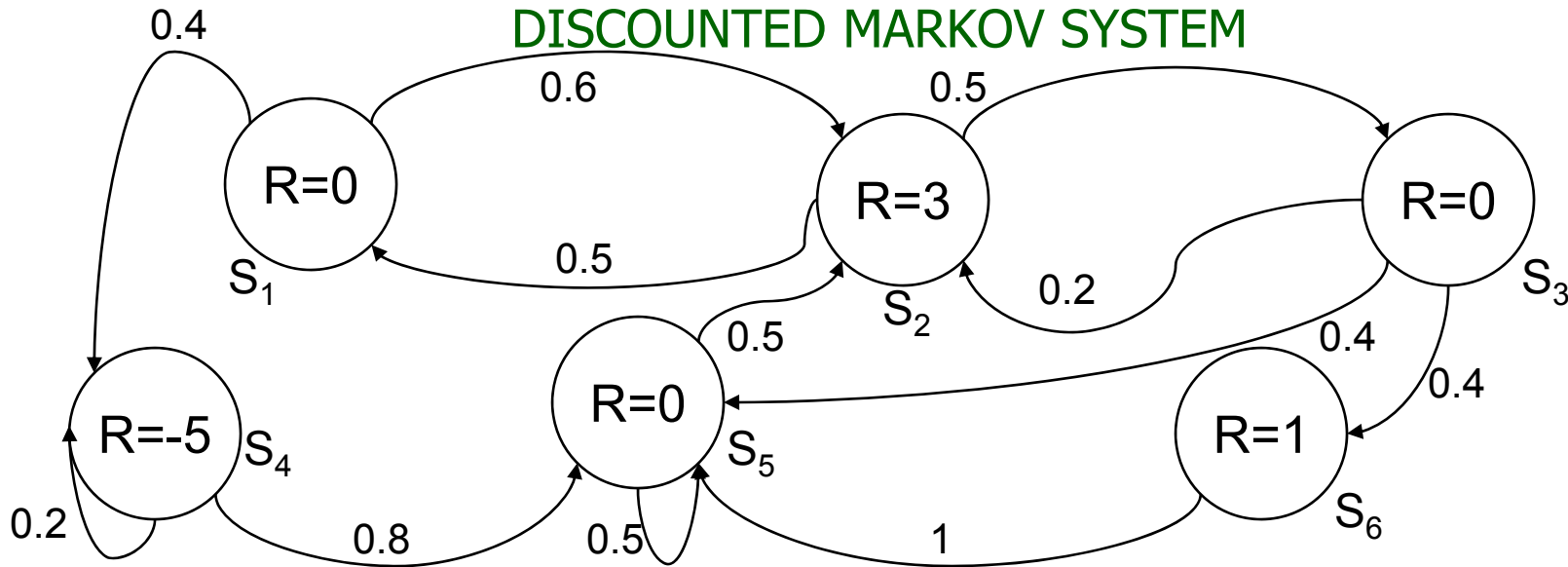


Reinforcement Learning

Predicting Delayed Rewards IN A DISCOUNTED MARKOV SYSTEM



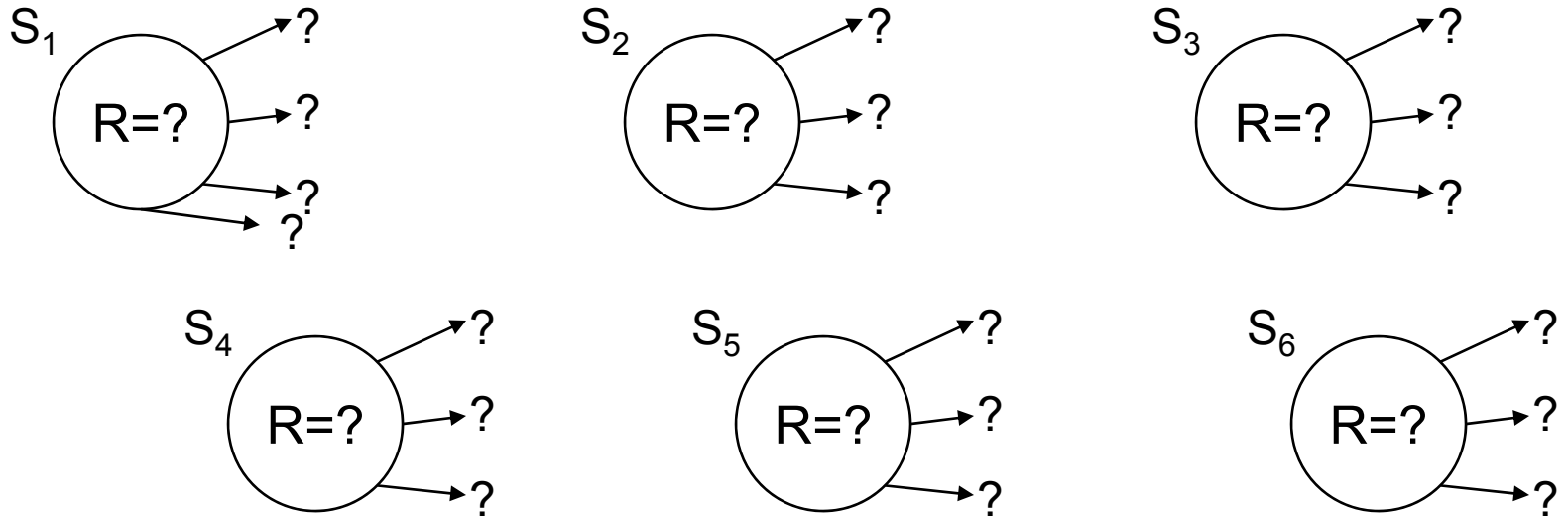
Prob(next state = S_5 | this state = S_4) = 0.8 etc...

What is expected sum of future rewards (discounted) ?

$$E \left[\left(\sum_{t=0}^{\infty} \gamma^t R(S[t]) \right) \mid S[0] = S \right]$$

Just Solve It! We use standard Markov System Theory

Learning Delayed Rewards...



All you can see is a series of states and rewards:

$S_1(R=0) \rightarrow S_2(R=0) \rightarrow S_3(R=4) \rightarrow S_2(R=0) \rightarrow S_4(R=0) \rightarrow S_5(R=0)$

Task: Based on this sequence, estimate $J^*(S_1), J^*(S_2), \dots, J^*(S_6)$

Idea 1: Supervised Learning

Assume $\gamma=1/2$

$S_1(R=0) \rightarrow S_2(R=0) \rightarrow S_3(R=4) \rightarrow S_2(R=0) \rightarrow S_4(R=0) \rightarrow S_5(R=0)$

At $t=1$ we were in state S_1 and eventually got a long term discounted reward of $0 + \gamma 0 + \gamma^2 4 + \gamma^3 0 + \gamma^4 0 \dots = 1$

At $t=2$ in state S_2 $l_{tdr} = 2$

At $t=5$ in state S_4 $l_{tdr} = 0$

At $t=3$ in state S_3 $l_{tdr} = 4$

At $t=6$ in state S_5 $l_{tdr} = 0$

At $t=4$ in state S_2 $l_{tdr} = 0$

State	Observations of LTDR	Mean LTDR	
S_1	1	1	$= J^{est}(S_1)$
S_2	2, 0	1	$= J^{est}(S_2)$
S_3	4	4	$= J^{est}(S_3)$
S_4	0	0	$= J^{est}(S_4)$
S_5	0	0	$= J^{est}(S_5)$

Supervised Learning ALG

- Watch a trajectory

$$S[0] r[0] S[1] r[1] \cdots S[T]r[T]$$

- For $t=0, 1, \dots, T$, compute $J[t] = \sum_{i=0}^{\infty} \gamma^i r[t+i]$

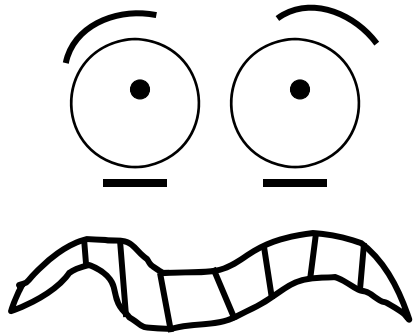
- Compute $J^{est}(S_i) = \left(\begin{array}{c} \text{mean value of } J[t] \\ \text{among all transitions beginning} \\ \text{in state } S_i \text{ on the trajectory} \end{array} \right)$

Let $\text{MATCHES}(S_i) = \{t \mid S[t] = S_i\}$, then define

$$J^{est}(S_i) = \frac{\sum_{t \in \text{MATCHES}(S_i)} J[t]}{|\text{MATCHES}(S_i)|}$$

- You're done!

Supervised Learning ALG



You may be worried about edge effects for some of the final transitions. With large trajectories these are negligible.

Online Supervised Learning

Initialize: $\text{Count}[S_i] = 0 \quad \forall S_i$
 $\text{SumJ}[S_i] = 0 \quad \forall S_i$
 $\text{Eligibility}[S_i] = 0 \quad \forall S_i$

Observe:

When we experience S_i with reward r
do this:

$\forall j \quad \text{Elig}[S_j] \leftarrow \gamma \text{Elig}[S_j]$
 $\text{Elig}[S_i] \leftarrow \text{Elig}[S_i] + 1$
 $\forall j \quad \text{SumJ}[S_j] \leftarrow \text{SumJ}[S_j] + r \times \text{Elig}[S_j]$
 $\text{Count}[S_i] \leftarrow \text{Count}[S_i] + 1$

Then at any time,

$J^{\text{est}}(S_j) = \text{SumJ}[S_j] / \text{Count}[S_j]$

Online Supervised Learning Economics

Given N states $S_1 \dots S_N$, OSL needs $O(N)$ memory.
Each update needs $O(N)$ work since we must update all
Elig[] array elements

Idea: Be sparse and only update/process Elig[]
elements with values $> \xi$ for tiny ξ

There are only $\log\left(\frac{1}{\xi}\right) / \log\left(\frac{1}{\gamma}\right)$
such elements

Easy to prove:

$$\text{As } T \rightarrow \infty, J^{est}(S_i) \rightarrow J^*(S_i) \quad \forall S_i$$

Online Supervised Learning



Let's grab OSL off the street, bundle it into a black van, take it to a bunker and interrogate it under 600 Watt lights.

$$S_1(r=0) \rightarrow S_2(r=0) \rightarrow S_3(r=4) \rightarrow S_2(r=0) \rightarrow S_4(r=0) \rightarrow S_5(r=0)$$

State	Observations of LTDR	$\hat{J}(S_i)$
S_1	1	1
S_2	2, 0	1
S_3	4	4
S_4	0	0
S_5	0	0

There's something a little suspicious about this (efficiency-wise)

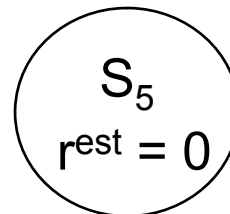
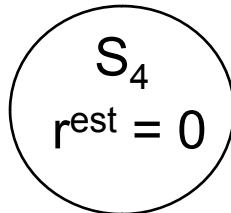
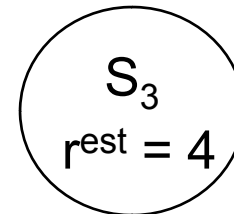
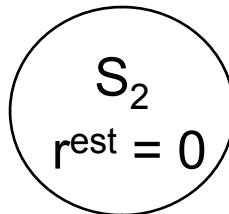
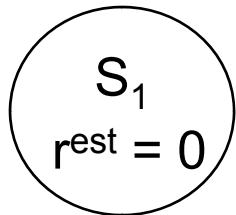
Certainty-Equivalent (CE) Learning

Idea: Use your data to estimate the underlying Markov system, instead of trying to estimate J directly.

$S_1(r=0) \rightarrow S_2(r=0) \rightarrow S_3(r=4) \rightarrow S_2(r=0) \rightarrow S_4(r=0) \rightarrow S_5(r=0)$

Estimated Markov System:

You draw in the transitions + probs



What're the estimated J values?

C.E. Method for Markov Systems

Initialize:

$$\left. \begin{array}{l} \text{Count}[S_i] = 0 \\ \text{SumR}[S_i] = 0 \\ \text{Trans}[S_i, S_j] = 0 \end{array} \right\} \begin{array}{l} \forall S_i \quad \# \text{Times visited } S_i \\ \forall S_i \quad \text{Sum of rewards from } S_i \\ \forall S_j \quad \# \text{Times transitioned from } S_i \rightarrow S_j \end{array}$$

When we are in state S_i , and we receive reward r , and we move to S_j ...

$$\text{Count}[S_i] \leftarrow \text{Count}[S_i] + 1$$

$$\text{SumR}[S_i] \leftarrow \text{SumR}[S_i] + r$$

$$\text{Trans}[S_i, S_j] \leftarrow \text{Trans}[S_i, S_j] + 1$$

Then at any time

$$r^{\text{est}}(S_j) = \text{SumR}[S_j] / \text{Count}[S_j]$$

$$P_{ij}^{\text{est}} = \text{Estimated Prob}(\text{next} = S_j \mid \text{this} = S_i)$$

$$= \text{Trans}[S_i, S_j] / \text{Count}[S_i]$$

C.E. for Markov Systems (continued) ...

So at any time we have

$$r^{est}(S_j) \text{ and } P^{est}(\text{next}=S_j \mid \text{this}=S_i)$$

$$\forall S_i S_j \quad = P_{ij}^{est}$$

So at any time we can solve the set of linear equations

$$J^{est}(S_i) = r^{est}(S_i) + \gamma \sum_{S_j} P^{est}(S_j \mid S_i) J^{est}(S_j)$$

[In vector notation,

$$\mathbf{J}^{est} = \mathbf{r}^{est} + \gamma \mathbf{P}^{est} \mathbf{J}$$

$$\Rightarrow \mathbf{J}^{est} = (\mathbf{I} - \gamma \mathbf{P}^{est})^{-1} \mathbf{r}^{est}$$

where \mathbf{J}^{est} \mathbf{r}^{est} are vectors of length N

\mathbf{P}^{est} is an $N \times N$ matrix

$N = \# \text{ states}$]

C.E. Online Economics

Memory: $O(N^2)$

Time to update counters: $O(1)$

Time to re-evaluate J^{est}

- $O(N^3)$ if use matrix inversion
- $O(N^2k_{\text{CRIT}})$ if use value iteration and we need k_{CRIT} iterations to converge
- $O(Nk_{\text{CRIT}})$ if use value iteration, and k_{CRIT} to converge, and M.S. is **Sparse** (i.e. mean # successors is constant)

Certainty Equivalent Learning

COMPLAINT

Memory use could be $O(N^2)$!

And time per update could be $O(Nk_{\text{CRIT}})$ up to $O(N^3)$!

Too expensive for some people.

Prioritized sweeping will help, (see later), but first let's review a very **inexpensive** approach

Why this obsession with online-ness?

I really care about supplying up-to-date J^{est} estimates all the time.

Can you guess why?

If not, all will be revealed in good time...

Less Time: More Data Limited Backups

- Do previous C.E. algorithm.
- At each time timestep we observe $S_i(r) \rightarrow S_j$ and update $\text{Count}[S_i]$, $\text{SumR}[S_i]$, $\text{Trans}[S_i, S_j]$
- And thus also update estimates

$$r_i^{est} \text{ and } P_{ij}^{est} \quad \forall_j \in \text{outcomes}(S_i)$$

But instead of re-solving for J^{est} , do **much less** work.
Just do one “backup” of $J^{est}[S_i]$

$$J^{est}[S_i] \leftarrow r_i^{est} + \gamma \sum_j P_{ij}^{est} J^{est}[S_j]$$

“One Backup C.E.” Economics

Space : $O(N^2)$

NO IMPROVEMENT
THERE!

Time to update statistics : $O(1)$

Time to update J^{est} : $O(1)$ 

- ❖ **Good News:** Much cheaper per transition
- ❖ **Good News:** Contraction Mapping proof (modified) promises convergence to optimal
- ❖ **Bad News:** Wastes data

Prioritized Sweeping

[Moore + Atkeson, '93]

Tries to be almost as data-efficient as full CE but not much more expensive than “One Backup” CE.

On every transition, some number (β) of states may have a backup applied. Which ones?

- The most “deserving”
- We keep a priority queue of which states have the biggest potential for changing their $J^{\text{est}}(S_j)$ value

Where Are We?

Trying to do online J^{est} prediction from streams of transitions

	Space	J^{est} Update Cost
Supervised Learning	$O(N_s)$	$O\left(\frac{1}{\log(1/\gamma)}\right)$
Full C.E. Learning	$O(N_{so})$	$O(N_{so}N_s)$ $O(N_{so}k_{\text{CRIT}})$
One Backup C.E. Learning	$O(N_{so})$	$O(1)$
Prioritized Sweeping	$O(N_{so})$	$O(1)$

Data Efficiency:



N_{so} = # state-outcomes (number of arrows on the M.S. diagram)

N_s = # states

What Next ?
Sample Backups !!!

Temporal Difference Learning

[Sutton 1988]

Only maintain a J^{est} array...
nothing else

So you've got

$J^{\text{est}}(S_1)$ $J^{\text{est}}(S_2)$, ... $J^{\text{est}}(S_N)$

and you observe

$S_i \xrightarrow{r} S_j$

what should you do?

A transition from i that receives
an immediate reward of r and
jumps to j

Can You Guess ?

TD Learning



We update $J^{est}(S_i)$

We nudge it to be closer to expected future rewards

$$J^{est}(S_i) \leftarrow (1 - \alpha)J^{est}(S_i) + \underbrace{\alpha \left[\begin{array}{c} \text{Expected future} \\ \text{rewards} \end{array} \right]}_{\text{WEIGHTED SUM}}$$
$$= (1 - \alpha)J^{est}(S_i) + \alpha \left[r + \gamma J^{est}(S_j) \right]$$

α is called a “learning rate” parameter. (See “ η ” in the neural lecture)

Decaying Learning Rate

[Dayan 1991ish] showed that for **General TD** learning of a Markov System (not just our simple model) that if you use update rule

$$J^{est}(S_i) \leftarrow \alpha_t [r_i + \gamma J^{est}(S_j)] + (1 - \alpha_t) J^{est}(S_i)$$

then, as number of observations goes to infinity $J^{est}(S_i) \rightarrow J^*(S_i) \forall i$

PROVIDED

- All states visited ∞ ly often

- $\sum_{t=1}^{\infty} \alpha_t = \infty$

- $\sum_{t=1}^{\infty} \alpha_t^2 < \infty$

This means

$$\forall k. \exists T. \sum_{t=1}^T \alpha_t > k$$

This means

$$\exists k. \forall T. \sum_{t=1}^T \alpha_t^2 < k$$

Decaying Learning Rate

This Works: $\alpha_t = 1/t$

This Doesn't: $\alpha_t = \alpha_0$

This Works: $\alpha_t = \beta/(\beta+t)$ [e.g. $\beta=1000$]

This Doesn't: $\alpha_t = \beta\alpha_{t-1}$ ($\beta < 1$)

A Fancier TD...

Write $S[t]$ = state at time t

Suppose $\alpha = 1/4$ $\gamma = 1/2$

Assume $J^{\text{est}}(S_{23})=0$ $J^{\text{est}}(S_{17})=0$ $J^{\text{est}}(S_{44})=16$

Assume $t = 405$ and $S[t] = S_{23}$

Observe $S_{23} \xrightarrow{(r=0)} S_{17}$ with reward 0

Now $t = 406$, $S[t] = S_{17}$, $S[t-1] = S_{23}$

$J^{\text{est}}(S_{23})=$, $J^{\text{est}}(S_{17})=$, $J^{\text{est}}(S_{44})=$

Observe $S_{17} \xrightarrow{(r=0)} S_{44}$

Now $t = 407$, $S[t] = S_{44}$

$J^{\text{est}}(S_{23})=$, $J^{\text{est}}(S_{17})=$, $J^{\text{est}}(S_{44})=$

INSIGHT: $J^{\text{est}}(S_{23})$ might think

What about me !!!

TD(λ) Comments







TD($\lambda=0$) is the original TD

TD($\lambda=1$) is almost the same as supervised learning (except it uses a learning rate instead of explicit counts)

TD($\lambda=0.7$) is often empirically the best performer

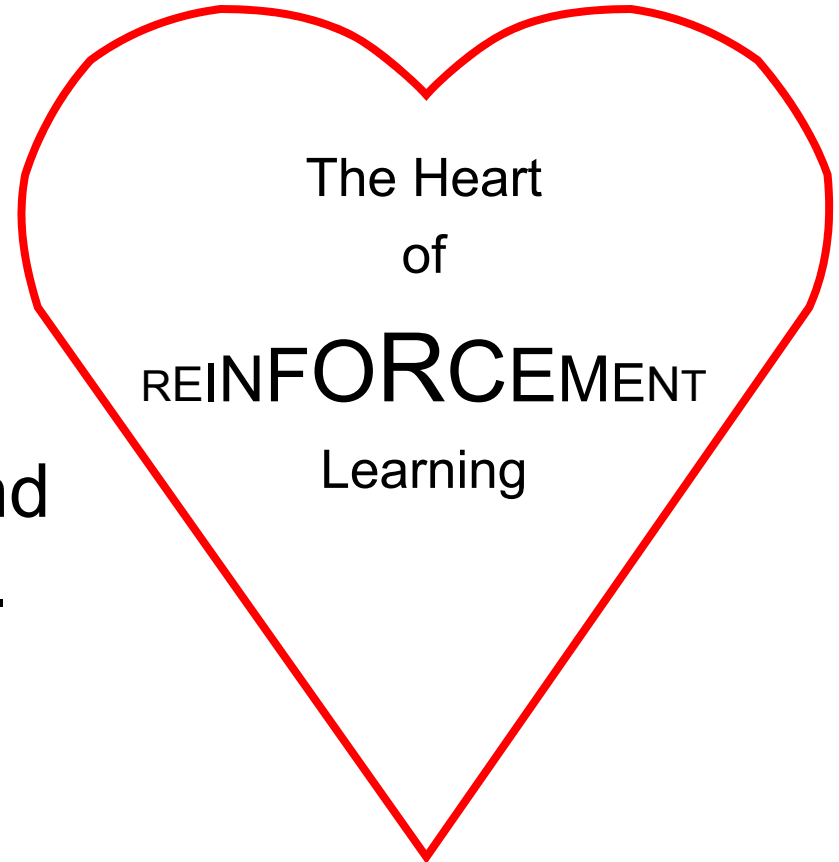
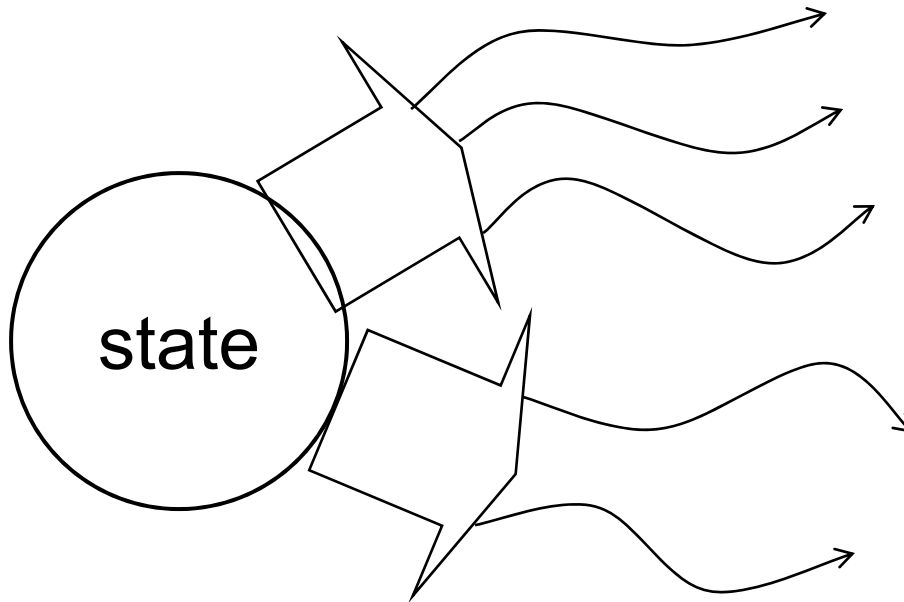
- Dayan's proof holds for all $0 \leq \lambda \leq 1$
- Updates can be made more computationally efficient with “eligibility” traces (similar to O.S.L.)
- Question:
 - ❖ Can you invent a problem that would make TD(0) look bad and TD(1) look good?
 - ❖ How about TD(0) look good & TD(1) bad??

Learning M.S. Summary

		Space	J Update Cost	Data Efficiency
MODEL-BASED	Supervised Learning	$O(N_s)$	$O\left(\frac{1}{\log \frac{1}{\gamma}}\right)$	
	Full C.E. Learning	$O(N_{so})$	$O(N_{so}N_s)$ $O(N_{so}k_{CRIT})$	
	One Backup C.E. Learning	$O(N_{so})$	$O(1)$	
	Prioritized Sweeping	$O(N_{so})$	$O(1)$	
MODEL FREE	TD(0)	$O(N_s)$	$O(1)$	
	TD(λ), $0 < \lambda \leq 1$	$O(N_s)$	$O\left(\frac{1}{\log \frac{1}{\gamma\lambda}}\right)$	

Learning Policies for MDPs

See previous lecture
slides for definition of and
computation with MDPs.



The task:

World: You are in state 34.

Your immediate reward is 3. You have 3 actions.

Robot: I'll take action 2.

World: You are in state 77.

Your immediate reward is -7. You have 2 actions.

Robot: I'll take action 1.

World: You're in state 34 (again).

Your immediate reward is 3. You have 3 actions.

The Markov property means once you've selected an action the P.D.F. of your next state is the same as the last time you tried the action in this state.

The “Credit Assignment” Problem

I'm in state 43,	reward = 0,	action = 2
“ “ “ 39,	“ = 0,	“ = 4
“ “ “ 22,	“ = 0,	“ = 1
“ “ “ 21,	“ = 0,	“ = 1
“ “ “ 21,	“ = 0,	“ = 1
“ “ “ 13,	“ = 0,	“ = 2
“ “ “ 54,	“ = 0,	“ = 2
“ “ “ 26,	“ = 100,	




Yippee! I got to a state with a big reward! But which of my actions along the way actually helped me get there??

This is the **Credit Assignment** problem.

It makes **Supervised Learning** approaches (e.g. **Boxes** [Michie & Chambers]) very, very slow.

Using the **MDP** assumption helps avoid this problem.

MDP Policy Learning

	Space	Update Cost	Data Efficiency
Full C.E. Learning	$O(N_{sA_0})$	$O(N_{sA_0} k_{\text{CRIT}})$	
One Backup C.E. Learning	$O(N_{sA_0})$	$O(N_{A_0})$	
Prioritized Sweeping	$O(N_{sA_0})$	$O(\beta N_{A_0})$	

- We'll think about **Model-Free** in a moment...
- The **C.E.** methods are very similar to the **MS** case, except now do value-iteration-for-**MDP** backups

$$J^{est}(s_i) = \max_a \left[r_i^{est} + \gamma \sum_{s_j \in \text{SUCCS}(s_i)} P^{est}(s_j | s_i, a) J^{est}(s_j) \right]$$

Choosing Actions

We're in state S_i

We can estimate r_i^{est}

“ “ “ $P^{est}(\text{next} = S_j \mid \text{this} = S_i, \text{action } a)$

“ “ “ $J^{est}(\text{next} = S_j)$

So what action should we choose ?

$$\text{IDEA 1: } a = \arg \max_{a'} \left[r_i + \gamma \sum_j P^{est}(S_j | S_i, a') J^{est}(S_j) \right]$$

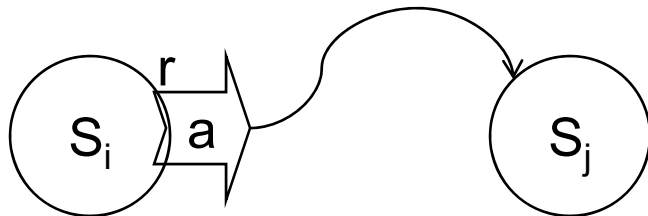
IDEA 2: $a = \text{random}$

- Any problems with these ideas?
- Any other suggestions?
- Could we be optimal?

Model-Free R.L.

Why not use T.D. ?

Observe



update

$$J^{est}(S_i) \leftarrow \alpha(r_i + \gamma J^{est}(S_j)) + (1 - \alpha)J^{est}(S_i)$$

What's wrong with this?

Q-Learning: Model-Free R.L.

[Watkins, 1988]

Define

$Q^*(S_i, a)$ = Expected sum of discounted future rewards if I start in state S_i , if I then take action a , and if I'm subsequently optimal

Questions:

Define $Q^*(S_i, a)$ in terms of J^*

Define $J^*(S_i)$ in terms of Q^*

Q-Learning Update

Note that

$$Q^*(S, a) = r_i + \gamma \sum_{S_j \in \text{SUCCS}(S_i)} P(S_j | S_i, \alpha) \max_{a'} Q^*(S_j, a')$$

In Q-learning we maintain a table of Q^{est} values instead of J^{est} values...

When you see S_i $\xrightarrow[\text{action } a]{\text{reward}}$ S_j do...

$$Q^{\text{est}}(S_i, a) \leftarrow \alpha \left[r_i + \gamma \max_{a'} Q^{\text{est}}(S_j, a') \right] + (1 - \alpha) Q^{\text{est}}(S_i, a)$$

This is even cleverer than it looks: the Q^{est} values are not biased by any particular exploration policy. It avoids the **Credit Assignment** problem.

Q-Learning: Choosing Actions

Same issues as for CE choosing actions

- Don't always be greedy, so don't always choose: $\arg \max_a Q(s_i, a)$
- Don't always be random (otherwise it will take a long time to reach somewhere exciting)

- Boltzmann exploration [Watkins]

$$\text{Prob}(\text{choose action } a) \propto \exp\left(-\frac{Q^{est}(s, a)}{K_t}\right)$$

- Optimism in the face of uncertainty [Sutton '90, Kaelbling '90]
 - Initialize Q-values optimistically high to encourage exploration
 - Or take into account how often each s,a pair has been tried

Q-Learning Comments

- [Watkins] proved that Q-learning will eventually converge to an optimal policy.
- Empirically it is cute
- Empirically it is very slow
- Why not do $Q(\lambda)$?
 - Would not make much sense [reintroduce the credit assignment problem]
 - Some people (e.g. Peng & Williams) have tried to work their way around this.

If we had time...

- Value function approximation
 - Use a Neural Net to represent J^{est} [e.g. Tesauro]
 - Use a Neural Net to represent Q^{est} [e.g. Crites]
 - Use a decision tree
 - ...with Q-learning [Chapman + Kaelbling '91]
 - ...with C.E. learning [Moore '91]
 - ...How to split up space?
 - Significance test on Q values [Chapman + Kaelbling]
 - Execution accuracy monitoring [Moore '91]
 - Game Theory [Moore + Atkeson '95]
 - New influence/variance criteria [Munos '99]

If we had time...

- R.L. Theory
 - Counterexamples [Boyan + Moore], [Baird]
 - Value Function Approximators with Averaging will converge to something [Gordon]
 - Neural Nets can fail [Baird]
 - Neural Nets with **Residual Gradient** updates will converge to something
 - Linear approximators for TD learning will converge to something useful [Tsitsiklis + Van Roy]