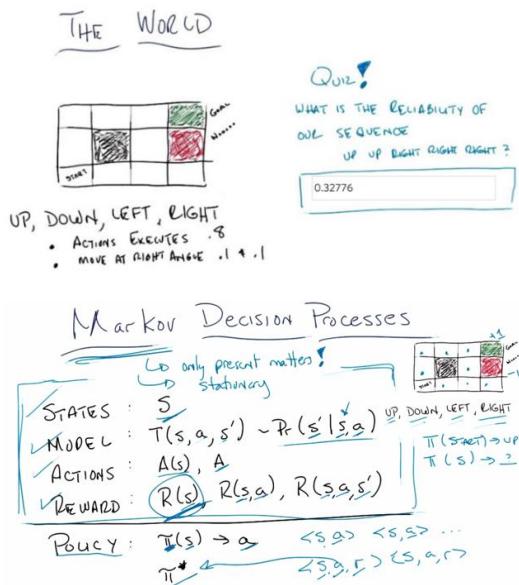


Week 1: Decision Making and Reinforcement Learning



States: set of tokens which describe every state that one could be in,

Model/Transition Function: probability that you'll end up transition to s' given you are in state s and taking action a . This is 'physics' or rules of the world (not changeable).

Action: $A(s)$, things you can do in a given state, a function of state, rule you can play; or A , a set actions not depends regardless of state.

Reward: scalar value that you get for being in a state. Usefulness of entering a state/and taking an action/and ending up into s'

S , T , A and R define the problem, policy is the solution,

Policy: the solution, a function that takes up a state and tells an action that you'll take $\langle s, a \rangle$ while a is the correct action you want to take to maximize the reward

Optimal policy: π^ , optimized, maximize long term expected reward. (Note, from $\langle s, a, r \rangle$, you know s , a and then know the r , based on rewards, find π^*)

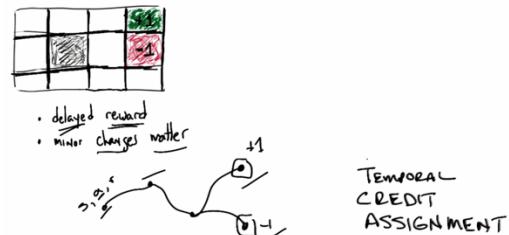
Recall that in Reinforcement Learning: $y = f(x)$ with z , given x and z , determine function f which can generate y . In Markov process, we know s and r , and we need to learn π to determine a . So, $s \rightarrow x$, $r \rightarrow z$, $a \rightarrow y$, $\pi \rightarrow f$

**Markovian property **:

--> Only the present matters; the transition state only depends on the current state

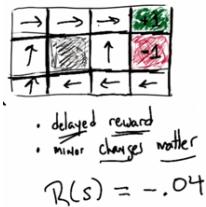
--> Transition model is stationary, rules doesn't change

MDPs: More About REWARDS



Delayed reward: in each state, we need figure out what action should we take to get the ultimate best rewards at the end.

MDPs: More About REWARDS



** Minor changes matter: By assigning small negative rewards in each state (except the absorbing state), the learning agent is encouraged to take one of the available actions to leave the current state and pursuit a better outcome.

If $R(s)$ is large positive, it encourages you to stay in the game. If $R(s)$ is large negative (< -1.6294), it encourages you to leave the world ASAP.

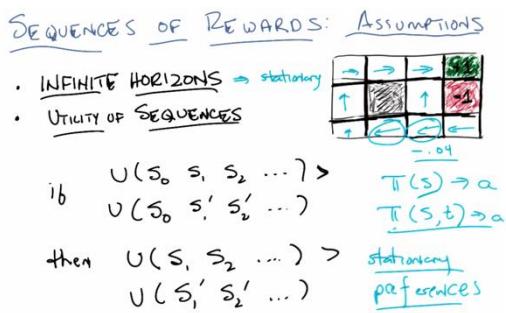
Reward is important to define learning behavior.

Rewards depends on domain knowledge.

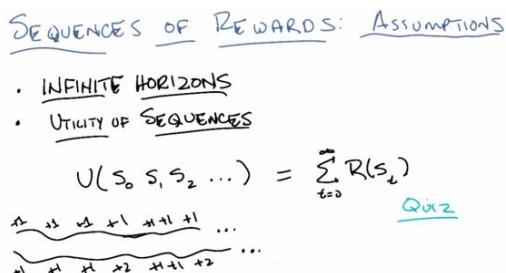
Changing the reward will change the goal of agent.



- Rationale: when the rewards are positive, the goal becomes never leaving the game. When the rewards are negative, the goal becomes leaving the game ASAP.



- Infinite horizons: if the time step is finite, the policy might change: $\pi(s, t) \rightarrow a$, and the Markovian property "π is stationary" will not stand anymore. So, we need to assume we have infinite horizons.



- Utility of Sequences: **stationary preferences**, Utility of all the states is the sum of the rewards of all states.

Since there are infinite number of states, the top sequence and the bottom sequence are essentially the same utility ∞ since the reward are all positive if we use the equation above.

So a new equation is introduced below by introducing a discount factor γ :

SEQUENCES OF REWARDS: Assumptions

INFINITE HORIZONS

UTILITY OF SEQUENCES

$$U(s_0, s_1, s_2, \dots) = \sum_{t=0}^{\infty} R(s_t)$$

discarded \Rightarrow geometric
infinite \Rightarrow finite

$$= \sum_{t=0}^{\infty} \gamma^t R(s_t) \quad 0 \leq \gamma < 1$$

$$\leq \sum_{t=0}^{\infty} \gamma^t R_{\max} = \frac{R_{\max}}{1-\gamma}$$

The equation introduced Discounted factor of Rewards and made the right side of the equation a geometric series so that the utility of sequences is bounded at $R_{\max}/(1-\gamma)$, where $0 \leq \gamma < 1$. This way, the infinite horizons is bounded by a finite number.

SEQUENCES OF REWARDS: Assumptions

$$(\sum_{t=0}^{\infty} \gamma^t) R_{\max}$$

$$x = (\gamma^0 + \gamma^1 + \gamma^2 + \dots)$$

$$\gamma^0 + \gamma^1 + \gamma^2 + \dots$$

$$x = \gamma^0 + \gamma^1$$

$$x - \gamma x = \gamma^0$$

$$x(1 - \gamma) = \gamma^0$$

$$x = 1 / 1 - \gamma \cdot R_{\max}$$

GEOMETRY EAST

POLICIES

$$\pi^* = \operatorname{argmax}_{\pi} E \left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \mid \pi \right]$$

$$R(s) \stackrel{\text{immediate}}{=} U(s) \stackrel{\text{long term}}{=} E \left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \mid \pi, s_0 = s \right]$$

delayed

Policy is A mapping from state to action

Optimal policy π^* : the policy allow us getting the largest expected discounted rewards when we follow the policy π^*

Utility of a state $U(s)$ is long term expected discounted rewards from **this point on**; It's delayed reward.

While Reward is immediate satisfaction.

$$\pi^*(s) = \operatorname{argmax}_a \sum_{s'} T(s, a, s') U(s')$$

$$U(s) = R(s) + \gamma \max_a \sum_{s'} T(s, a, s') U(s')$$

Bellman Equation

$\pi^*(s)$ is the optimal policy given current state. **
 $U(s)^*$ is the utility of the state if follow π^* . Thus

we get the utility of the current state is reward of current state $R(s)$ plus the discounted utility from this point on. This is the Bellman Equation.

POLICIES : FINDING POLICIES

$$U(s) = R(s) + \gamma \max_a \sum_{s'} T(s, a, s') U(s')$$

n equations in n unknowns

- start w/ arbitrary utilities
update utilities based on neighbors
repeat until converge

$$\hat{U}(s) = R(s) + \gamma \max_a \sum_{s'} T(s, a, s') \hat{U}(s')$$

The Bellman equation has n equations and n unknowns. Because the equations are non-linear (because the max operation), the equations are not easy to solve.

Here list the way to solve the problem: by setting an **arbitrary utility** and update utility based on neighbors, and repeat until converge. This method is called **value iteration**. Neighbors are any state that we can reach.

The reason that VI can work is the reward of each state is true and the utility is discounted.

POLICIES : FINDING POLICIES

$$U(s) = R(s) + \gamma \max_a \sum_{s'} T(s, a, s') U(s')$$

$$U_{t+1}(s) = R(s) + \gamma \max_a \sum_{s'} T(s, a, s') U_t(s')$$

QUIT IF

$\gamma = 1/2$	$R(s) = -0.14$	$U_0(s) = 0$
$-0.04 + \gamma [0.2 + 0.8]$	$U_1(s) = .36$	
$-0.04 + \gamma [0.36 + 0.8]$	$U_2(s) = .3716$	

- In this quiz, assuming the initial U of all states but the terminal states is 0. Compute the U at first and second value iteration.

POLICIES : FINDING POLICIES

Policy Iteration

- start with $\pi_0 \leftarrow \text{guess}$
- evaluate: given π_t calculate $U_t = U^{T_{\pi_t}}$
- improve: $\pi_{t+1} = \arg \max_a \sum_{s'} T(s, a, s') U_t(s')$

$$U_{t+1}(s) = R(s) + \gamma \sum_{s'} T(s, \pi_t(s), s') U_t(s')$$

n equations in n unknowns

By starting with a guessed policy, the Bellman equation don't have to do Max when calculating the expected utility. So there are n linear equation with n unknowns, and it's fairly easy to get the estimation of utility.

The Bellman Equation

$$V(s) = \max_a (R(s, a) + \gamma \sum_{s'} T(s, a, s') V(s'))$$

$$V(s) = \max_a (R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_a (R(s', a') + \gamma \sum_{s''} T(s', a', s'') \dots))$$

$V(s)$ is infinite sequence with sub V . By regrouping the Value function, we can get Q function as:

The Bellman Equations

Value $V(s) = \max_a (R(s, a) + \gamma \sum_{s'} T(s, a, s') V(s'))$

Utility $Q(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_a Q(s', a')$

$$V(s) = \max_a (R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_a (R(s', a') + \gamma \sum_{s''} T(s', a', s'') \dots))$$

And C function as:

$C(s, a) = \gamma \sum_{s'} T(s, a, s') \max_{a'} (R(s', a') + C(s', a'))$

The Third Bellman Equation

$Q(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} (R(s', a') + \gamma \sum_{s''} T(s', a', s'') \max_{a''} (R(s'', a'') + \dots))$

$C(s) = \gamma \sum_{s'} T(s, a, s') \max_{a'} (R(s', a') + C(s', a'))$

$C(s, a) = \gamma \sum_{s'} T(s, a, s') \max_{a'} (R(s', a') + C(s', a'))$

$C(s, a, r) = \gamma \sum_{s'} T(s, a, s') \max_{a'} (r + C(s', a', r'))$

$C(s, a, s') = \gamma \sum_{s''} T(s, a, s') \max_{a''} (R(s'', a'') + C(s', a', s''))$

continuation

Quiz 5: the correct form of C function

- In the first equation, we do not know which action a to use to transition into state s' , so we cannot compute the transition function.
- Correct answer. In this formula, all the variables are properly defined, and there is a meaningful relationship between the continuation function and the reward function.

- Notice that there is no meaningful dependency of the max operator on the parameter r (which is already passed in to the continuation function). Also, the variable r' is never defined.
- In this equation, the state s'' is not defined, so it does not make sense to pass it into the recursive call of the continuation function.

What's the difference between V form and Q form and C form of Bellman equation?

- In value function, V is connected by transition T and reward r. we must know transition and the reward of current state S. V is updated with respect to a state.
- In quality function, we can rely on experience to update Q, this is useful in reinforcement learning when we do not know T and R ahead of time. Q-value is updated with respect to an action.

	V	Q	C
V	$V(s) = V(s)$	$V(s) = \max_a Q(s, a)$	$V(s) = \max_a (R(s, a) + \gamma \sum_{s'} T(s, a, s') V(s'))$
Q	$Q(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') V(s')$	$Q(s, a) = Q(s, a)$	$Q(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') Q(s', a)$
C	$C(s, a) = \max_{s'} T(s, a, s') V(s')$	$C(s, a) = \max_{s'} T(s, a, s') \max_a Q(s', a)$	$C(s, a) = C(s, a)$

The Bellman Equations

Value $V(s) = \max_a (R(s, a) + \gamma \sum_{s'} T(s, a, s') V(s'))$
Quality $Q(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_a Q(s', a')$

$$V(s) = \max_{a_1} (R(s_1, a_1) + \gamma \sum_{s_2} T(s_1, a_1, s_2) V(s_2)) = \max_{a_1} (R(s_1, a_1) + \gamma \sum_{s_2} T(s_1, a_1, s_2) \max_{a_2} (R(s_2, a_2) + \gamma \sum_{s_3} T(s_2, a_2, s_3) V(s_3)) \dots$$

$$\square C(s, a) = \gamma \sum_{s'} T(s, a, s') \max_{a'} (R(s', a') + C(s', a'))$$

We don't need reward function when we convert Q to C or V. We can convert C to V and q without knowing the transition function.

What Have We Learned?

- Relearned MDPs (Scooby Doo MDP)
 S, A, T, R, \dots
 $\stackrel{R(s)}{\parallel}$
 $R(s, a)$
- Q functions (relation to V and C)
- C function (continuations) *notation matters*
- Bellman Equations *representation*

Week 2: Intro to Burlap

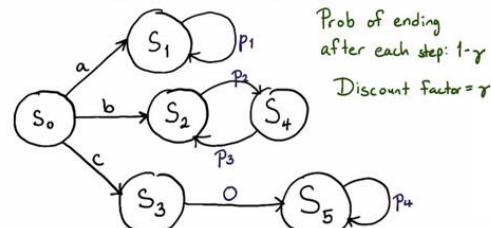
In this introduction lecture to Burlap, I learned how to use BURLAP to setup a MDP.

Steps includes:

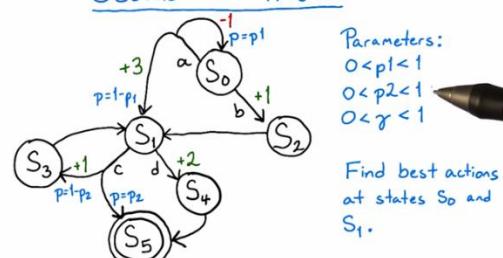
- setup domain (using GraphDefinedDomain)
- setup nodes (states)
- setup actions
- initialize state
- setup reward function, termination function, hashFactory
- run Value Iteration
- And get the optimal first action.

See the sample code below and read [Burlap documentation](#) to understand the code.

First MDP Problem



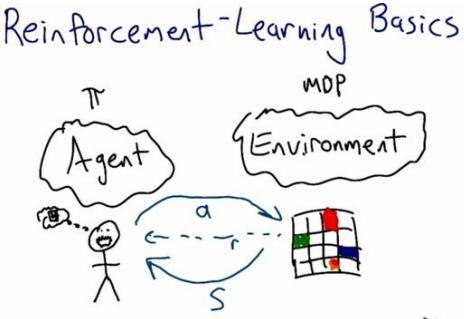
Second MDP Problem



See code: First MDP Problem

The solution of this second problem will not be posted here since it's the quiz of the lecture.

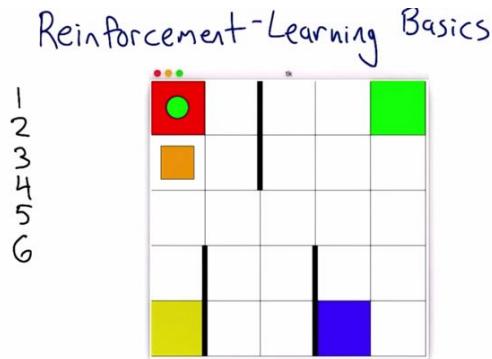
Week 3: Reinforcement Learning Basics



In RL, environment is only available to agent as perceived states (s), the agent can interact with the environment by taking action (a) and the environment gives a reward (r) as feedback to tell the agent if the $\langle s, a \rangle$ pair are good or not. The computation is calculated in the agent's head.

The difference of RL and MDP is that in MDP, environment is totally available to the agent, while in RL, Environment is only available through the agent's perception.

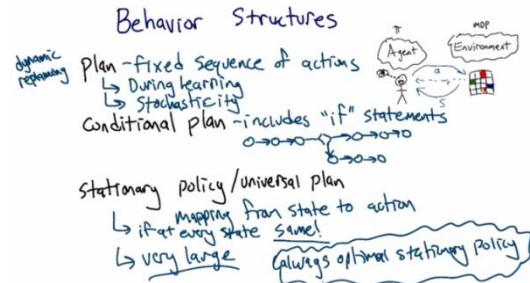
Demo of RL



The small orange square represents the agent, and it can perform 6 actions. The world is the grid with some colored squares and a green dot. The goal what's the game and what are the actions.

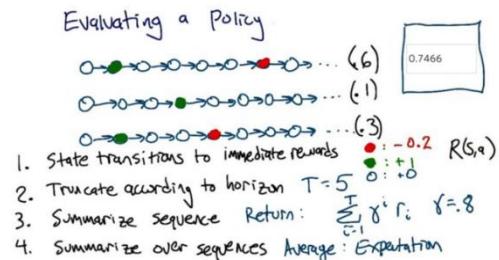
Behavioral structure

The goal is to generate learning algorithm



- **Plan** is a set of fixed actions. Plan won't work during learning or when the environment is only partially known or stochastic.
- **Conditional Plan** includes "if" statements
- **Stationary policy (or Universal Plan)** are mapping from state to action. It can handle stochastic very well but it is very large. **There always is an optimal stationary policy.**

Evaluating a policy



Quiz: evaluating a policy

- The numbers in the parentheses are probabilities of choosing the sequence. $R(s,a)$ is reward function. **Return** is discounted rewards.
- $0.8^1 * 0.6 + 0.8^3 * 0.1 + (0.8^1 + 0.8^4 / (-0.2)) * 0.3 = 0.746624$
- This is the expected value of the policy based on the assumption that we index the states in each sequence from left to right starting at 0 on the far left. Interpreting $T=5$ to mean either truncating at the fifth circle or after the fifth transition (i.e. at the sixth circle) in each row gives the same result.

Evaluating a Learner

Evaluating a Learner

- Value of returned Policy
- Computational complexity (time)
- Experience complexity (time)
↳ how much data it needs

MIMIC



Recap

What Have we Learned ?

- What RL is: agent \rightarrow environment maximizing rewards (hard!)
- What a policy is
↳ plans, policies
- Evaluating plays
- Evaluating learners efficiency, effectiveness

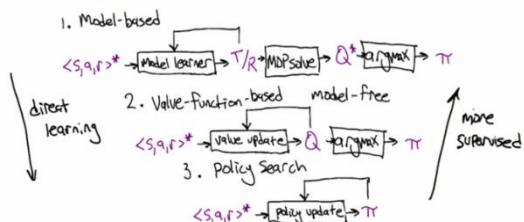
Week 4: Temporal Difference Learning

Temporal Difference Learning

Read Sutton (1988)

- Read Sutton, Read Sutton, Read Sutton. Because the final project was based on it!

RL Context

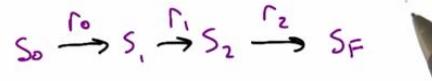


- Model based
- Model free
- Policy search
- Form 1 --> 3: more direct learning
- From 3 --> 1 more supervised

TD-lambda

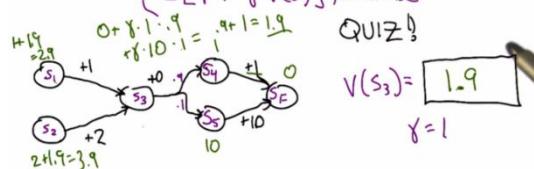
$$TD(\lambda)$$

Learning to predict over time.



Example

$$\text{Learn } V(s) = \begin{cases} 0, & \text{if } s = S_F \\ E[r + \gamma V(s')] & \text{otherwise} \end{cases}$$



- In this case the model is known, the calculation is easy.

Estimating From Data

1. $s_1 \xrightarrow{a_1} s_3 \xrightarrow{r_1} s_F \xrightarrow{a_2} r_2$
 2. $s_1 \xrightarrow{a_1} s_3 \xrightarrow{r_1} s_F \xrightarrow{a_2} r_2$
 3. $s_1 \xrightarrow{a_1} s_3 \xrightarrow{r_1} s_F \xrightarrow{a_2} r_2$
 4. $s_1 \xrightarrow{a_1} s_3 \xrightarrow{r_1} s_F \xrightarrow{a_2} r_2$
 5. $s_2 \xrightarrow{a_1} s_3 \xrightarrow{r_1} s_F \xrightarrow{a_2} r_2$
- QUIZ: What is an appropriate estimate for $V(s_1)$ after 3 episodes? 4?

$$\frac{15}{3} = 5$$

$$\frac{17}{4} = 4.25$$

- Remember from the previous lecture, we need to get value from each episode and average over them.

Computing Estimates Incrementally

$$V_{t-1}(s_t) \quad R_t(s_t) \quad V_t(s_t)$$

$$\begin{aligned} V_t(s_t) &= \frac{(T-1)V_{t-1}(s_t) + R_t(s_t)}{T} \\ &= \frac{T-1}{T} V_{t-1}(s_t) + \frac{1}{T} R_t(s_t) \\ &= V_{t-1}(s_t) + \alpha_T (R_t(s_t) - V_{t-1}(s_t)) \end{aligned}$$

where $\alpha_T = \frac{1}{T}$ error

- The rewrite makes the formula looks a lot like neuro-net learning. And alpha is introduced.

Properties of Learning Rates	
$V_T(s) = V_{T-1}(s) + \alpha_T (R_T(s) - V_{T-1}(s))$	① $\sum \alpha_T = \infty$
$\lim_{T \rightarrow \infty} V_T(s) = V(s)$	② $\sum \alpha_T^2 < \infty$
$\sum \alpha_T$	$\sum \alpha_T^2 / T$
$\alpha_T = 1/T$	harmonic
$\alpha_T = 1/T^{0.5}$	
$\alpha_T = 1/T^{0.2}$	
$\alpha_T = 1/T^{0.1}$	
$\alpha_T = 1/T^{-0.1}$	
$\alpha_T = 1/T^{-0.2}$	
$\alpha_T = 1/T^{-0.5}$	
$\alpha_T = 1/T^{-0.9}$	
$\alpha_T = 1/T^{-1}$	
$\alpha_T = 1/T^{-1.1}$	
$\alpha_T = 1/T^{-1.5}$	
$\alpha_T = 1/T^{-2}$	

Quiz 2

Quiz 2: alpha will make learning converge (tips: if $\alpha_i > 1$, $1/(T)^{i-1}$ will be bounded)

TD(1) Rule

Episode T
For all s_t , $e(s_t) = 0$ at start of episode, $V_T(s) = V_{T-1}(s)$
After $s_{t+1} \xrightarrow{r_t} s_t$: (step t)
 $e(s_{t+1}) = e(s_{t+1}) + 1$
For all s_t
 $V_T(s) = V_{T-1}(s) + \alpha_T (r_t + \gamma V_{T-1}(s_t) - V_{T-1}(s_{t+1}))e(s)$
 $e(s) = \gamma e(s)$

TD(1) rule

TD(1) Example
 $s_1 \xrightarrow{r_1} s_2 \xrightarrow{r_2} s_3 \xrightarrow{r_3} s_f$
 $e = \gamma^2 r_2 + r_1$
 $\Delta V(s_1) \propto (r_1 + \gamma r_2 + \gamma^2 r_3 + \gamma^3 V_{T-1}(s_f) - V_{T-1}(s_1))$
 $\Delta V(s_2) \propto (r_2 + \gamma r_3 + \gamma^2 V_{T-1}(s_f) - V_{T-1}(s_2))$
 $\Delta V(s_3) \propto (r_3 + \gamma V_{T-1}(s_f) - V_{T-1}(s_3))$
 r_1, r_2, r_3, r_f
CLAIM:
TD(1) is the same as outcome-based updates (if no repeated states)
Now with extra learning!

TD(1) with and without repeated states

- When no repeated states, the TD(1) is the same as outcome-based updates (which is see all the rewards in each state and update weights).
- When there is repeated states, extra learning happens.

Why TD(1) is "wrong"?

$s_1 \xrightarrow{r_1} s_3 \xrightarrow{r_2} s_4 \xrightarrow{r_3} s_f$
 $s_1 \xrightarrow{r_1} s_3 \xrightarrow{r_2} s_5 \xrightarrow{r_3} s_f$
 $s_1 \xrightarrow{r_1} s_3 \xrightarrow{r_2} s_4 \xrightarrow{r_3} s_f$
 $s_1 \xrightarrow{r_1} s_3 \xrightarrow{r_2} s_4 \xrightarrow{r_3} s_f$
 $s_1 \xrightarrow{r_1} s_3 \xrightarrow{r_2} s_5 \xrightarrow{r_3} s_f$
 $\xrightarrow{r_4} s_1 \xrightarrow{r_1} s_3 \xrightarrow{r_2} s_5 \xrightarrow{r_3} s_f$

→ average across all encounters

$\bar{x} = 1$

QUIZ $V(s_2)$?
Outcome-based (TD(1)) estimate?
+12
Maximum likelihood estimate?
+6.6

Why TD(1) is "Wrong"

- In case of TD(1) rule, $V(s_2)$ can be estimated by average episodes. We only see $V(s_2)$ once and the value is 12. Then $V(s_2) = 12$
- In case of Maximum likelihood estimates, we have to kind of learn the transition from data. E.g. for the first 5 episodes, we saw $s_3 \rightarrow s_4$ 3 times and $s_3 \rightarrow s_5$ 2 times. So the transition probability can be extracted from data as 0.6 and 0.4 respectively.

TD(0) Rule

Finds ML estimate (if data repeated infinitely often)

$$V_T(s_{t+1}) = V_T(s_{t+1}) + \alpha_T (r_t + \gamma V_T(s_t) - V_T(s_{t+1}))$$

$$V_T(s_{t+1}) = \underset{s_t}{E} [r_t + \gamma V_T(s_t)]$$

$$V_T(s_{t+1}) = E [r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots]$$

TD(0) Rule

- First of all, if we have infinite data, TD(1) will also do the right thing.
- When we have finite data, we can repeatedly infinitely sample the data to figure out all the ML. This is what TD(0) do.

TD(0) Rule

Episode T
For all s_t , $e(s_t) = 0$ at start of episode, $V_T(s) = V_{T-1}(s)$
After $s_{t+1} \xrightarrow{r_t} s_t$: (step t)
 $e(s_{t+1}) = e(s_{t+1}) + 1$
For all s_t
 $V_T(s) = V_T(s) + \alpha_T (r_t + \gamma V_{T-1}(s_t) - V_{T-1}(s_{t+1}))e(s)$
 $e(s) = \gamma e(s)$

TD(0)
Episode T
For all s_t , $e(s_t) = 0$ at start of episode, $V_T(s) = V_{T-1}(s)$
After $s_{t+1} \xrightarrow{r_t} s_t$: (step t)
 $e(s_{t+1}) = e(s_{t+1}) + 1$
For all s_t
 $V_T(s) = V_T(s) + \alpha_T (r_t + \gamma V_{T-1}(s_t) - V_{T-1}(s_{t+1}))e(s)$
 $e(s) = \gamma e(s)$

TD(0)
Episode T
For all s_t , $e(s_t) = 0$ at start of episode, $V_T(s) = V_{T-1}(s)$
After $s_{t+1} \xrightarrow{r_t} s_t$: (step t)
 $e(s_{t+1}) = e(s_{t+1}) + 1$
For all s_t
 $V_T(s) = V_T(s) + \alpha_T (r_t + \gamma V_{T-1}(s_t) - V_{T-1}(s_{t+1}))e(s)$
 $e(s) = \gamma e(s)$

Connecting TD(0) and TD(1)

K-Step Estimators

K-Step Estimators

$$\begin{aligned}
 E_1 \quad V(S_t) &= V(S_t) + \alpha_T(r_{t+1} + \gamma V(S_{t+1}) - V(S_t)) \quad TD(0) \\
 E_2 \quad V(S_t) &= V(S_t) + \alpha_T(r_{t+1} + \gamma r_{t+2} + \gamma^2 V(S_{t+2}) - V(S_t)) \\
 E_3 \quad V(S_t) &= V(S_t) + \alpha_T(r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 V(S_{t+3}) - V(S_t)) \\
 E_K \quad V(S_t) &= V(S_t) + \alpha_T(r_{t+1} + \dots + \gamma^{K-1} r_{t+K} + \gamma^K V(S_{t+K}) - V(S_t)) \\
 E_\infty \quad V(S_t) &= V(S_t) + \alpha_T(r_{t+1} + \dots + \gamma^{K-1} r_{t+K} + \dots + \gamma^\infty V(S_t)) \quad TD(\infty)
 \end{aligned}$$

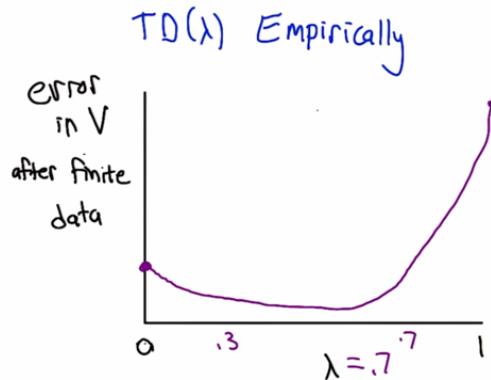
K-Step Estimators

- E1 is one-step estimator (one-step look up) TD(0)
- E2 is two-step estimator, and Ek is k-step lookup.
- When K goes to infinity, we got TD(1)

$$\begin{array}{lll}
 \text{K-Step Estimators and TD}(\lambda) & & \\
 \begin{matrix} \lambda=0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{matrix} & \begin{matrix} \lambda=1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{matrix} & \begin{matrix} \lambda=\frac{1}{2} \\ 0 \\ \frac{1}{2} \\ 0 \\ \frac{1}{2} \\ 0 \\ \frac{1}{2} \\ 0 \\ \frac{1}{2} \end{matrix} \\
 \begin{matrix} TD(0) \\ TD(1) \\ E_1 \\ E_2 \\ E_3 \\ E_K \\ \dots \\ E_\infty \end{matrix} & &
 \end{array}
 \begin{array}{l}
 E_1 \quad 1-\lambda \\
 E_2 \quad \lambda(1-\lambda) \\
 E_3 \quad \lambda^2(1-\lambda) \\
 E_K \quad \lambda^{K-1}(1-\lambda) \\
 E_\infty \quad \lambda^\infty = \lambda \cdot \frac{1}{1-\lambda} = 1
 \end{array}$$

K-step Estimators and TD-lambda

TD-lambda can be seen as weighted combination of K-step estimators. The weight factor are $\lambda^k(1-\lambda)$.



Why use TD-lambda?

The best performed lambda is typically not TD(0), but some λ in between 0 and 1.

What Have We Learned?

- Temporal difference learning - value estimates
- Context : RL
 - model-based
 - value-based
 - policy-based
- Incremental estimates, outcome based
- Learning rate properties $\sum \alpha_t = \infty, \sum \alpha_t^2 < \infty$
- TD(1) - inefficient - high variance
- TD(0) - ML estimate, cleaner estimates
- TD(λ) - interpolates, cleaner estimates

Week 5: Convergence: TD with control

Bellman Equations

$$\begin{aligned}
 \text{no actions: } V(S) &= R(S) + \gamma \sum_{S'} T(S, S') V(S') \\
 \langle S_{t-1}, r_t, S_t \rangle & \Rightarrow V(S_t) = V_{t-1}(S_{t-1}) + \alpha_t(r_t + \gamma V_{t-1}(S_t) - V_{t-1}(S_{t-1})) \\
 V_{t-1}(S) &= V_{t-1}(S), \text{ otherwise} \quad TD(0) !
 \end{aligned}$$

Learning without control:

- no actions by learner. The value function is that value equals reward $R(s)$ plus discounted expected value of the next state. ($T(s, s')$ is transition function).
- The new estimated value function for the state the learner just left is TD(0).

Bellman Equations with Actions

$$\begin{aligned}
 \text{actions: } Q(S, a) &= R(S, a) + \gamma \sum_{S'} T(S, a, S') \max_a Q(S', a) \\
 \langle S_{t-1}, a_{t-1}, r_t, S_t \rangle & \Rightarrow Q_{t-1}(S_{t-1}, a_{t-1}) = Q_{t-1}(S_{t-1}, a_{t-1}) + \alpha_t(r_t + \gamma \max_a Q_{t-1}(S_t, a) - Q_{t-1}(S_{t-1}, a_{t-1})) \\
 Q_t(S, a) &= Q_{t-1}(S, a), \text{ otherwise} \quad TD(0) \\
 \text{Approximations: } & \text{① If we knew the model, synchronously update} \\
 Q_{t-1}(S, a) &= R(S, a) + \gamma \sum_{S'} T(S, a, S') \max_a Q_{t-1}(S', a) \\
 & \text{② If we knew } Q^*, \text{ Sampling asynchronously update} \\
 Q_t(S_{t-1}, a_{t-1}) &= Q_{t-1}(S_{t-1}, a_{t-1}) + \alpha_t(r_t + \gamma \max_a Q^*(S_t, a) - Q_{t-1}(S_{t-1}, a_{t-1}))
 \end{aligned}$$

Learning with control/actions.

Q function is used to estimate the value at current state, take a action then get a reward and ended up in a new state $\langle S_{t-1}, a_{t-1}, r_t, S_t \rangle$. The Q updating rule takes care of two approximations: 1) if model is known, it can be

used to update Q; 2) if Q^* is known, it can also be used to update Q.

And both will converge.

Bellman operator

Let B be an operator, or mapping from value functions to value functions

$$[BQ](s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q(s', a')$$

Quiz $Q^* = BQ^*$ is another way of writing

$$\begin{aligned} Q_0 &= 0, & \text{Bellman Equation!} \\ Q_t = BQ_{t-1} &\text{ is another way of writing} & \\ & \quad \text{Value Iteration!!!} \end{aligned}$$

Quiz 1: Bellman Operator

See value iteration [here](#)

Contraction Mappings

B is an operator $\|Q\|_\infty = \max_{s, a} |Q(s, a)|$

If, for all F, G and some $0 \leq \gamma < 1$

$$\|BF - BG\|_\infty \leq \gamma \|F - G\|_\infty$$

then B is a contraction mapping.

Contraction mapping def: If applying the B operator makes the distance between two functions smaller than the the distance between the original functions.

Contraction Mapping Quiz

In the space of real numbers, which of the following are contraction mappings?

- $B(x) = \frac{x}{2}$
- $B(x) = x+1$
- $B(x) = x-1$
- $B(x) = (x+100) \cdot 0.9 = x \cdot 0.9 + 90$

- Practically, if B is multiplying a number in $[0, 1]$, B is contraction mapping.
- Take the first option for example: $\|B(x) - B(y)\| \leq \gamma \|x - y\|$. So $B(x) = x/2$ is a contraction mapping.

Contraction Properties

If B is a contraction mapping,

① $B^* = BF^*$ has a solution and it is unique

② $F_t = BF_{t-1} \Rightarrow F_t \rightarrow F^*$
(value iteration converges)

$$\|F_t - F^*\|_\infty = \|BF_{t-1} - BF^*\|_\infty \leq \gamma \|F_{t-1} - F^*\|_\infty$$

Contraction Properties

- ① and ② are true, so that $BF_{t-1} = F_t$, F_t will converge at F^* through value iteration.
- If there is two fix point, G^* and F^* , Putting them into the B operator will not change the distance of them because both of them are fixed, and this violates the definition of B .

Bellman Operator Contracts

$$[BQ](s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q(s', a')$$

Given Q_1, Q_2

$$\begin{aligned} \|BQ_1 - BQ_2\|_\infty &= \max_{s, a} |[BQ_1](s, a) - [BQ_2](s, a)| = \\ &= \max_{s, a} |(R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_1(s', a')) - (R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_2(s', a'))| \\ &= \max_{s, a} |\gamma \sum_{s'} T(s, a, s') (\max_{a'} Q_1(s', a') - \max_{a'} Q_2(s', a'))| \\ &\leq \gamma \|Q_1 - Q_2\|_\infty \end{aligned}$$

- Applying B operator to Bellman equation and unpacking $\|BQ_1 - BQ_2\|_\infty$.
- Not combining the $\max a'$ because the a' are not the same in Q_1 and Q_2

Bellman Operator Contracts

Given Q_1, Q_2

$$\begin{aligned} \|BQ_1 - BQ_2\|_\infty &= \max_{s, a} |[BQ_1](s, a) - [BQ_2](s, a)| = \\ &= \max_{s, a} |(R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_1(s', a')) - (R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_2(s', a'))| \\ &= \max_{s, a} |\gamma \sum_{s'} T(s, a, s') (\max_{a'} Q_1(s', a') - \max_{a'} Q_2(s', a'))| \\ &\leq \max_{s, a} |\max_{a'} Q_1(s', a') - \max_{a'} Q_2(s', a')| \quad \text{if} \\ &\leq \gamma \max_{s, a} |Q_1(s, a) - Q_2(s, a)| = \gamma \|Q_1 - Q_2\|_\infty \end{aligned}$$

- Stop here, first, will go back. That call this the unfinished proof.

Max is a non-expansion

for all $f, g \quad |\max_a f(a) - \max_a g(a)| \leq \max_a |f(a) - g(a)|$

QUIZ $\boxed{\quad} \leq \boxed{\quad}$

	a_1	a_2	a_3	a_4	a_5	a_6	a_7	a_8
f	7	9	1	0	6	3	-2	
g	5	3	8	1	8	5	5	

Quiz 3

Max is a non-expansion

for all $f, g \quad |\max_a f(a) - \max_a g(a)| \leq \max_a |f(a) - g(a)|$

Proof: WLOG $\max_a f(a) \geq \max_a g(a)$

$$|\max_a f(a) - \max_a g(a)| = \max_a f(a) - \max_a g(a)$$

$$= f(a_1) - g(a_1) \quad \begin{matrix} f(a_1) \\ g(a_1) \end{matrix} \quad \left\{ \begin{matrix} a_1 = \arg \max_a f(a) \\ a_2 = \arg \max_a g(a) \end{matrix} \right.$$

$$\leq f(a_1) - g(a_1)$$

$$= |f(a_1) - g(a_1)|$$

$$\leq \max_a |f(a) - g(a)|$$

- Max is non-expansion: $||\max f(a) - \max g(a)|| \leq \max |f(a) - g(a)|$. Using this

Convergence Define $\alpha_t(s,a) = 0$ if $s_t \neq s, a \neq a$

Theorem: Let B be a contraction mapping and $Q^* = BQ^*$ be its fixed point. Let Q_0 be a γ^t function and define $Q_{t+1} = [B + Q_t]Q_t$. Then, $Q_t \rightarrow Q^*$ if:

- For all $U_1, U_2, S, a : |([B + U_1]Q^*)(S, a) - ([B + U_2]Q^*)(S, a)| \leq (\alpha_t(S, a)) |U_1(S, a) - U_2(S, a)|$
- For all $Q, U, S, a : |([B + U]Q^*)(S, a) - ([B + U]Q)(S, a)| \leq \alpha_t(S, a) |Q^*(S, a) - Q(S, a)|$
- $\sum_t \alpha_t = \infty$, $\sum_t \alpha_t^2 < \infty$ $\leq \alpha_t(S, a) |Q^*(S, a) - Q(S, a)|$

Statement of theorem

The three properties need to be true for Q_t to converge to Q^* . And they are true

So, Q-learning converges.

Convergence Theorem Explained

$$[B + Q]w(s, a) = Q(s, a) + d_t(s, a) (r_t + \gamma \max_{a'} w(s', a'))$$

$$1. \text{ For all } U_1, U_2, S, a : |([B + U_1]Q^*)(S, a) - ([B + U_2]Q^*)(S, a)| \leq (\alpha_t(S, a)) |U_1(S, a) - U_2(S, a)|$$

- B_t is the operator we are going to update Q at t time step. $Q(s, a)$ is the Q function value of the state we just left and Q function w is the Q value of the state we just arrive. In

regular Q Learning update, Q and w are the same, here we separated them in the theorem.

- So, rule number one is saying if we know the Q^* and use the updating rule (in pink) to update the Q function, the expected value of the one-step look ahead (w) can be calculated and the stochasticity will be averaged out.

Convergence Theorem Explained

$$[B + Q]w(s, a) = Q(s, a) + d_t(s, a) (r_t + \gamma \max_{a'} w(s', a') - Q(s', a))$$

$$2. \text{ For all } Q, U, S, a : |([B + U]Q^*)(S, a) - ([B + U]Q)(S, a)| \leq \alpha_t(S, a) |Q^*(S, a) - Q(S, a)|$$

CONTRACTION

- If we hold the $Q(s, a)$ fixed and only varies the way we calculate the one-step look ahead, the distance between the Q^* and Q can only get closer with each update.
- The third condition is the learning rate condition. Which is needed for Bellman equation.

Generalized MDPs

$$Q^*(s, a) = R(s, a) + \gamma \sum_{s'} \sum_{a'} p(s'|s, a) Q^*(s', a')$$

$$\textcircled{1} g(s') = \sum_{s'} T(s, s', s) g(s') \quad \textcircled{2} f(s', a') = \max_{a'} f(s', a')$$

$$\textcircled{3} g(s') = \min_{s'} g(s') \quad \textcircled{4} f(s', a') = \max_{a'} f(s', a')$$

$$\textcircled{5} g(s') = \sum_{s'} T(s, s', s) g(s') \quad \textcircled{6} f(s', a') = \sum_{a'} p(s'|s, a) f(s', a')$$

$$\textcircled{7} g(s') = \sum_{s'} T(s, s', s) g(s') \quad \textcircled{8} f(s', a') = \min_{a'} \max_{s'} f(s', a')$$

regular MDP

□

□

□

Quiz 4: 1. Decision making on estimated value based on best next action (regular MDP). 2. The environment puts you in the worst possible state and you choose the next best action given the state (risk averse); 3. The state is the expected value but the action to take is based on ranking of the action (exploration-sensitive, min, max, mediocre); 4. (Zero-sum game)

- The take home messages are, these generalized MDPs all converges.

- So the correct answers are: (from top to bottom) worst possible decision, not-optimal decision, and two agent competing decision making (the zero - sum game).

What Have we Learned?

- convergence — Q-learning to Q^* , Value Iteration
 - generalized convergence theorem
 - contraction, non-expansion
 - generalized MDPs
- order statistics
fixed convex combinations
- coco-Q learning
when
converges

Recap

Generalized MDP can be seen as redefine fix point. Contraction might be something like "收敛" in Chinese. Q-learning converges to Q^* . Generalized Convergence theorem uses two Q-functions to prove the convergence of Bellman Equation.

Week 6: Advanced Algorithmic Analysis

Value iteration

More on VI

1. For some t^* polynomial in $|S|, |A|$, $R_{\max} = \max_{a,s} |R(s,a)|$, $\frac{1-\gamma}{1-\gamma}$, bits of precision, $T(s) = \arg \max_a Q_{t^*}(s, a)$ is optimal. Cramer's rule
2. If $|V_t(s) - V_{t+1}(s)| < \epsilon$ $\forall s$ $H \sim \frac{1}{1-\gamma}$
3. $\|B^K Q_1 - B^K Q_2\| \leq \gamma^K \|Q_1 - Q_2\|$.

- 1 tells us that VI converges in a reasonable amount of time t^* (finite). We know t^* exists, but we don't know when that will happen.
- 2 Gives us a bound of difference between the policy now and optimal policy. We can take this bound and test when it is a decent time to stop VI and take the current policy.
- Both 1 and 2 encourage us choosing a small γ . (γ tells us how further in the future we

should look). The effective horizon is $H \sim 1/(1-\gamma)$.

- 3 applying Bellman Operator K times to Q functions will shrink the distance between Q functions.

Linear Programming

- Only one way to solve MDP in polynomial time: solving Bellman equation through Linear Programming. A *linear programming* problem may be defined as the problem of maximizing or minimizing a linear function subject to linear constraints. (Definition extracted from this [PDF](#))
- Bellman equation has one part that is not linear, the *max* function. But it can be expressed by a series of linear function and a objective function *min*.

Linear Programming

$$\begin{aligned} \forall s, v_s &= \max_a (R(s,a) + \gamma \sum_{s'} T(s,a,s') v_{s'}) \\ \text{PRIMAL} \\ \forall s, v_s &= \min \sum_a v_a \\ & R(s,a) + \gamma \sum_{s'} T(s,a,s') v_{s'} \end{aligned}$$

- In the *primal*, the objective function is the minimum of sum of all the Vs;
- In linear programming, we can change constraints to variables and variables to constraints, and the resulting linear program is equivalent to the old one.

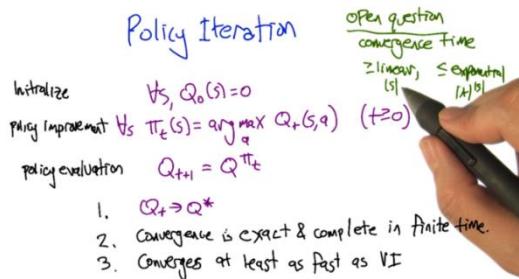
Linear Programming

$$\begin{aligned} \max_{q_{sa}} \sum_{s,a} q_{sa} R(s,a) & \quad \text{"Policy flow"} \\ \text{DUAL} \\ \text{s.t. } 1 + \gamma \sum_{s,a} q_{sa} T(s,a,s') &= \sum_a q_{sa} \cdot v_{s'} \\ v_{s'} & \geq 0 \end{aligned}$$

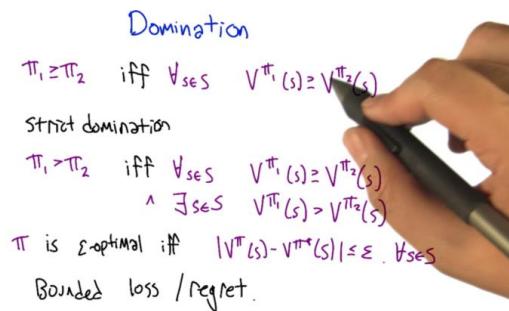
- Dual is a new linear program comes from the old primal version of linear program (没有推导过程)

- q_{sa} is "Policy flow", maximize the expected rewards of all states.
- For each possible next state, we wanted it to be true that the amount of policy flows going through the next state should be equal to the number of the state it has been visited.

Policy Iteration



- Initialize the first step Q value of all the states to be 0; improve the policy at time t; Apply the policy to calculate the Q value of t+1 step.
- Convergence time is an open question (but it is finite): \geq linear, \leq exponential



- For every state, if the value of it follows π_1 always equals or is larger than the value when it follows π_2 , we say π_1 **dominates** π_2 .
- If π_1 **dominates** π_2 and there exists states that $V^{\pi_1}(s) > V^{\pi_2}(s)$, we say π_1 **strictly dominates** π_2 .
- If for every state, the distance between the value following policy π and following the optimal policy π^* is no larger than ϵ , then π is ϵ -optimal.

Why Does Policy Iteration Work?

A tale of two operators for two policies:

π_1, π_2

$$[B_1 V](s) = R(s, \pi_1(s)) + \gamma \sum_{s'} T(s, \pi_1(s), s') V(s')$$

$$[B_2 V](s) = R(s, \pi_2(s)) + \gamma \sum_{s'} T(s, \pi_2(s), s') V(s')$$

- B_1 makes the update follows π_1 and B_2 makes the update follows π_2
- Applying B operator to two Values will not make them further apart, if the two values are the same, applying B will not make them closer.

B_2 is Monotonic

$$[B_2 V](s) = R(s, \pi_2(s)) + \gamma \sum_{s'} T(s, \pi_2(s), s') V(s')$$

$$[B_2 V](s) = R(s, \pi_2(s)) + \gamma \sum_{s'} T(s, \pi_2(s), s') V(s')$$

Theorem: If $V_1 \geq V_2$, then $B_2 V_1 \geq B_2 V_2$
Proof: $(B_2 V_1 - B_2 V_2)(s) = \gamma \sum_{s'} T(s, \pi_2(s), s') (V_1(s') - V_2(s')) \geq 0$

- Theorem: if V_1 dominates V_2 , applying B will keep the ordering: $BV_1 \geq BV_2$.

Another Property in Policy Iteration

$$\begin{aligned} \pi_1 &\rightarrow B_1 \rightarrow Q_1 = B_1 Q_1 \\ \pi_2 &\text{ greedy policy wrt } Q_1 \\ &\downarrow \\ B_2 & Q_1 \leq B_2 Q_1 \quad B_2 Q_1 \geq Q_1 \\ Q_1(s, a) &= R(s, a) + \gamma \sum_{s'} T(s, \pi_1(s), s') Q_1(s', \pi_1(s')) \\ &\leq R(s, a) + \gamma \sum_{s'} T(s, \pi_2(s), s') Q_1(s', \pi_2(s')) \end{aligned}$$

monotonicity property
Value improvement
defn of $\pi_2, B_1, B_2, Q_1, Q_2$

- Q_1 is the fix point of B_1 ;
- π_2 is greedy policy with respect to Q_1
- B_2 is the Bellman operator of π_2
- Applying B_2 (the greedy function wrt Q_1) on Q_1 will add a bound to Q_1 , thus getting a better Q_1 . This is called **Value Improvement**.

Policy Iteration Proof

π_1 policy, B_1 operator, Q_1 value function wrt π_1 ,
 π_2 greedy wrt Q_1 , B_2 operator, Q_2 value function wrt π_2

$$\begin{aligned} & B_2 Q_1 \geq Q_1 \\ \text{KZO, } & B_2^{k+1} Q_1 \geq B_2^k Q_1 \\ \lim_{k \rightarrow \infty} & B_2^k Q_1 \geq Q_1 \\ & Q_2 \geq Q_1 \end{aligned}$$

Policy Iteration Proof

π_1 policy, B_1 operator, Q_1 value function wrt π_1 ,
 π_2 greedy wrt Q_1 , B_2 operator, Q_2 value function wrt π_2
non-deprovement

$$\begin{aligned} & B_2 Q_1 \geq Q_1 && \text{value improvement} \\ \text{KZO, } & B_2^{k+1} Q_1 \geq B_2^k Q_1 && \text{monotonicity} \\ \lim_{k \rightarrow \infty} & B_2^k Q_1 \geq Q_1 && a \geq b, b \geq c \Rightarrow a \geq c \\ & Q_2 \geq Q_1 && \text{fixed point} \end{aligned}$$

- Value improvement (or value non deprovement): for each state, value will improved until it could not get better anymore.
- Monotonicity: Value can only get better after each iteration.
- Transitivity: $a \geq b, b \geq c$, so $a \geq c$.
- Fixed point. If we apply B_2 over and over again on Q_1 , we will reach the fixed point of B_2 , which is Q_2 .

What Have We Learned?

- VI → converges in finite steps (greedy policy)
 - DILP → Primal / dual values - flow
 - PI → domination (state-by-state basis)
value non deprovement / value improvement
no local optima
monotonicity
- PROOFS DO NOT STINK

- In Value iteration, greedy policy converges in finite step. This does not necessarily mean value function will converge.

Week 7: messing with reward

Messing With Rewards



(Ng, Harada, Russell)

This week's tasks

- Watch Reward Shaping.
- Read Ng, Harada, Russell (1999) and* Asmuth, Littman, Zinkov (2008).*
- Office hours on Friday, October 2nd, from 4-5 pm(EST).
- Homework 6.

Why changing RF?

Why might we want to change the reward function for an MDP?

- ① Semantics
- ② Efficiency
 - speed
 - space
 - solvability

- Easier to solve / represent and similar to what it would have learned.
- Because we don't have one yet.

Given an MDP, RF can affect the behavior of the learner/agent so it ultimately specifies the behavior (or policy) we want for the MDP. So changing rewards can make the MDP easy to solve and represent

- Semantics:** what the agent are expected to do at each state;
- Efficiency:** speed (experience and/or computation needed), space (complexity), and solvability.

How to change RF without changing optimal policy.

How can we change the MDP
reward function without changing
the optimal policy? $\langle S, A, R, T, \gamma \rangle$

- multiply by positive constant
- shift by constant (add)
- non-linear potential-based



Given an MDP described by $\langle S, A, R, T, \gamma \rangle$, there are three ways to change R without changing optimal solution. (Note, if we know T , then it is not a RL problem anymore, so this part of lecture is for MDP not RL specifically).

- Multiply by a positive constant (non-zero 'cause multiply by 0 will erase all the reward information)
- shift by a constant
- non-linear potential-based

1. Multiply by a positive constant

Multiply by a (positive) scalar

$$Q(s,a) = R(s,a) + \gamma \sum_{s'} T(s,a,s') \max_{q'} Q(s',q')$$

$$R'(s,a) = c R(s,a) \quad (c > 0)$$

$$Q'(s,a) = \boxed{c Q(s,a)}$$

$$c Q(s,a) = c R(s,a) + \gamma \sum_{s'} T(s,a,s') \max_{q'} c Q(s',q')$$

- $Q(s,a)$ is the solution of Bellman function with the old RF $R(s,a)$.
- $R'(s,a)$ is a new RF with is the old RF multiplying by a constant.
- What's the new solution $Q'(s,a)$ with respect to the new RF $R'(s,a)$ and old $Q(s,a)$?

Here is how to solve the problem:

- $Q = R + \gamma R + \gamma^2 R + \dots + \gamma^\infty R$
- $Q' = R' + \gamma R' + \gamma^2 R' + \dots + \gamma^\infty R'$

- Replace R' with $c*R$,
- $$\begin{aligned} Q' &= (c*R) + \gamma(c*R) + \gamma^2(c*R) + \dots + \gamma^\infty(c*R) \\ &= c(R + \gamma R + \gamma^2 R + \dots + \gamma^\infty R) \\ &= c*Q \end{aligned}$$

2. Shift by a constant

Adding a Scalar

$$Q(s,a) = R(s,a) + \gamma \sum_{s'} T(s,a,s') \max_{q'} Q(s',q')$$

$$R'(s,a) = R(s,a) + c$$

$$Q'(s,a) = \boxed{\quad}$$

Adding a Scalar

$$\begin{aligned} Q(s,a) &= R(s,a) + \gamma \sum_{s'} T(s,a,s') \max_{q'} Q(s',q') \\ R'(s,a) &= R(s,a) + c \end{aligned}$$

$$Q'(s,a) = \boxed{Q(s,a) + \frac{c}{1-\gamma}}$$

$$\begin{aligned} Q(s,a) + \frac{c}{1-\gamma} &= R(s,a) + c + \gamma \sum_{s'} T(s,a,s') \left(\max_{q'} Q(s',q') \right) \frac{c}{1-\gamma} \\ \frac{c}{1-\gamma} &= c + \frac{c\gamma}{1-\gamma} \Leftrightarrow c = (1-\gamma)c + c\gamma \end{aligned}$$

Solution and proof of Quiz 2

- $Q = R + \gamma R + \gamma^2 R + \dots + \gamma^\infty R$
 - $Q' = R' + \gamma R' + \gamma^2 R' + \dots + \gamma^\infty R'$
 - Replace R' with $R+c$,
- $$\begin{aligned} Q' &= (R+c) + \gamma(R+c) + \gamma^2(R+c) + \dots + \gamma^\infty(R+c) \\ &= (R + \gamma R + \gamma^2 R + \dots + \gamma^\infty R) + (c + \gamma c + \gamma^2 c + \dots + \gamma^\infty c) \end{aligned}$$
- The first part is Q and the second part is geometric series. So,
- $$Q' = Q + c/(1-\gamma)$$

3. Nonlinear potential-based reward shaping

Potential-based Shaping

$$Q(s,a) = \sum_{s'} T(s,a,s') (R(s,a,s') + \gamma \max_{a'} Q(s',a'))$$

$$R'(s,a,s') = R(s,a) - \Psi(s) + \gamma \Psi(s')$$

$$Q'(s,a) =$$

Quiz 3: potential-based reward shaping

- $Q = R + \gamma R + \gamma^2 R + \dots + \gamma^\infty R$
- $Q' = R' + \gamma R' + \gamma^2 R' + \dots + \gamma^\infty R'$
- Replace R' with $R - \Psi(s) + \gamma \Psi(s')$,
 $Q' = (R - \Psi(s) + \gamma \Psi(s')) + \gamma(R - \Psi(s') + \gamma \Psi(s'')) + \dots + \gamma^\infty (R - \Psi(s^\infty) + \gamma \Psi(s^\infty))$

$$= (R + \gamma R + \gamma^2 R + \dots + \gamma^\infty R) + (-\Psi(s) + \gamma \Psi(s') + \gamma(-\Psi(s') + \gamma \Psi(s'')) + \gamma^2(-\Psi(s') + \gamma \Psi(s'')) + \dots + \gamma^\infty(-\Psi(s^\infty) + \gamma \Psi(s^\infty)))$$

- The first part is Q . In the second part, most of the elements are cancelling each other out and only has the very first and last elements left. So,
 $Q' = Q + (-\Psi(s) + \gamma^\infty \Psi(s^\infty))$
- Given γ is in $(0,1)$, so $\gamma^\infty = 0$. Then we have Q' :
 $Q' = Q - \Psi(s)$

Q-learning with potential

Updating the Q function with the potential based reward shaping,

- Q function will converge at $Q^*(s,a)$.
- We know that $Q(s,a) = Q^*(s,a) - \Psi(s)$. If we initialize $Q(s,a)$ with zero, then $Q^*(s,a) - \Psi(s) = Q^*(s,a) - \max_a Q^*(s,a) = 0$, that means a is optimal.
- So Q-learning with potential is like initializing Q at Q^*

Q-learning with Potentials

$$Q(s,a) := Q(s,a) + \alpha_r (r - \Psi(s) + \gamma \Psi(s') + \gamma \max_{a'} Q(s',a') - Q(s,a))$$

$$\begin{aligned} \Psi(s) &= \max_a Q^*(s,a) & Q(s,a) &= Q^*(s,a) - \Psi(s) \\ \text{we have} & & & - \max_a Q^*(s,a) \\ Q(s,a) &= 0 & \underline{Q(s,a)} &= 0 \quad \text{if } a \text{ optimal} \\ & & & \geq 0 \quad \text{otherwise} \end{aligned}$$

Q learning with potential \equiv Q learning with α -function initializat

What have we learned?

What Have We Learned?

- You can mess with rewards (not Texas)
 - ↳ modify rewards to shape agent experience without changing the optimal policy
- scale (positively) ? linear change w.r.t. the function
- shift (add a constant)
- potential functions ← speed things up. suboptimal
- Reward shaping and its dangers (positive loops)

- Potential functions is a way to speed up the process to solve MDP
- Reward shaping might have suboptimal positive loops which will never converge?

Week 8: Exploration

This week

- Watch Exploration.
- The readings were *Fong (1995) *and Li, Littman, Walsh (2008).

Topics

- bandits
- deterministic MDPs
- stochastic MDPs

Type	state transition	Stoch-astic	solution	
Bandits	x	✓	hoeffding bound to do stochastic decision making	• Maximum Confidence: ① Given current info, figure out the best expected action based on best confidence -> ② act, and get new info -> repeat ① ② will also choose b only
Deterministic MDPs	✓	x	Rmax to do sequential decision making	• Minimum Confidence: ① Given current info, figure out the best expected action based on smallest confidence -> ② act, and get new info -> repeat ① ② will also choose to alternating between a and b. But it will not favor the better arm.
Stochastic MDPs	✓	✓	both hoeffding bound and Rmax	

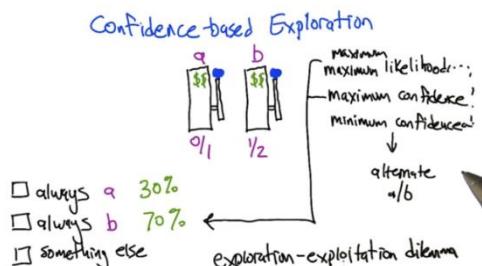
K-armed Bandit Problem

- Each slot has a probability to get certain payoff or no payoff
- Payoff or the associated probabilities are unknown.
- What's the best thing to do? **Exploration**.
- Exploration-exploitation dilemma: combine maximum likelihood and minimum confidence. The latter allows to try different actions and the former allows us to choose the best reward.

Confidence based exploration

!Quiz 1: K-Armed Bandits, given the observed data, which arm gives the highest expected payoff and which arm are most confident] (http://upload-images.jianshu.io/upload_images/118382-6578df35f8c00073.png?imageMogr2/auto-orient/strip%7CimageView2/2/w/1240)

- Arm d gives 1 payoff in the smallest number of tries
- Arm g gives most reliable output since it's been pulled the most and gives a good payoff.



- Maximum Likelihood: ① Given current info, figure out the best expected action based on best likelihood -> ② act, and get new info -> repeat ① ② will always choose b

Metrics for Bandits

Metrics for Bandits

1. identify optimal arm in the limit
2. maximize (discounted) expected reward
Gittins index $\text{GI} \stackrel{s}{\rightarrow} \frac{s}{d} \rightarrow \text{index}$
3. maximize expected reward over finite horizon
4. identify near-optimal arm (ε) with high probability $(1-\delta)$ in time $\tau(k, \varepsilon, \delta)$ (poly in $k, \delta, \frac{1}{\varepsilon}$)
5. nearly maximize reward (ε) with high probability $(1-\delta)$ in time $\tau'(k, \varepsilon, \delta)$ (poly in $k, \delta, \frac{1}{\varepsilon}$)
6. pull a non-optimal arm (ε) no more than $\tau''(k, \varepsilon, \delta)$ times with high probability $(1-\delta)$

- identify optimal arm in the limit: exploration
- maximize (discounted) expected rewards: exploration + exploitation, but computationally expensive
 - **Gittins index**: a real scalar value associated to the state of a stochastic process with a reward function and a probability of termination.
 - Gittins index works well with Bandits, should not generalize to other RL.
- Maximize expected reward in finite horizon - similar to value iteration
- Identify near-optimal arms (ε) with high probabilities $(1-\delta)$ in time $\tau(k, \varepsilon, \delta)$

(polynomial in k , $1/\epsilon$, $1/\delta$) - The "get there" goal is to find near optimal arms. Similar to PAC learning: probably approximately correct.

- Nearly Maximize reward (ϵ) with high probabilities ($1-\delta$) in time $\tau'(k, \epsilon, \delta)$ (polynomial in k , $1/\epsilon$, $1/\delta$): the "how to get there" goal is to accumulate as much more rewards as possible in finite time τ' .
- Pull a non-near optimal arm (ϵ) no more than $\tau''(k, \epsilon, \delta)$ times with high probabilities ($1-\delta$): The "control error" goal is to minimize mistakes.

Find Best Implies Few Mistakes

Find Best \Rightarrow Few mistakes

Algorithm that finds ϵ best with probability $1-\delta$ in $\tau(k, \epsilon, \delta)$ pulls. (A)

Want ϵ' close arm on $\tau'(k, \epsilon', \delta')$ pulls with probability $1-\delta'$.

$\epsilon = \epsilon'$, $\delta = \delta'$, run algorithm A and after $\tau(k, \epsilon, \delta)$ pulls we know ϵ' close arm. Pull it forever. $\tau'(k, \epsilon', \delta') \leq \tau(k, \epsilon, \delta)$

- τ' is the total number of mistakes when we find ϵ' close arm
- τ' is bounded to τ .

Few Mistakes Implies Do Well

Few mistakes \Rightarrow Do well

We have an algorithm that pulls ϵ -suboptimal arms at most $\tau(k, \epsilon, \delta)$ times with probability $1-\delta$. (B)

Want an algorithm C that, with high probability $1-\delta'$, has a per-step reward ϵ' close to optimal after $\tau'(k, \epsilon', \delta')$ steps.

$$\epsilon' = \frac{\epsilon}{2}$$

$$\tau'(k, \epsilon', \delta') = n$$

$$\tau(k, \epsilon, \delta) = m$$

$$m = \frac{\epsilon}{\epsilon'} = \frac{\epsilon}{\frac{\epsilon}{2}} = 2$$

$$n = \frac{m}{1-\delta} = \frac{2}{1-\delta}$$

$$\tau'(k, \epsilon', \delta') = n = \frac{2}{1-\delta}$$

- Is it true that the smaller the m is, the smaller the n is? not always.

- But it is true that the smaller the m is, the small the n could be.
- "Do well" means algorithm C can lead to ϵ' close to optimal arm where $\epsilon' < \epsilon$.

Do Well Implies Find Best

Do well \Rightarrow Find Best

We have an algorithm that gets within ϵ (per step) of optimal after $\tau(k, \epsilon, \delta)$ pulls with probability $1-\delta$.

Want an algorithm A which, after $\tau(k, \epsilon, \delta)$ pulls, returns an arm that is ϵ' optimal. up to ϵ' .

Run C: per step near optimal payoff.

$$\epsilon_i: \text{suboptimality of the plurality arm}$$

$$\epsilon_i \leq \epsilon' \quad \epsilon = \epsilon'/k \quad \epsilon_i \leq \epsilon/k$$

Putting It Together

Find Best \Rightarrow Few mistakes

Few mistakes \Rightarrow Do well

Do well \Rightarrow Find Best

Don't have to pick.



- Given algorithm of one goal, we can derive algorithms to reach the other two.
- So there is no need to pick one.
- Need to know that the three goals are not always solvable (so you might have to pick one?).

Hoeffding



Let X_1, \dots, X_n be iid with mean μ . Let $\hat{\mu}$ be our estimate of μ : $\sum_i X_i / n$. The following is a $100(1-\delta)\%$ confidence interval for μ :

$$\left[\hat{\mu} - \frac{z\delta}{\sqrt{n}}, \hat{\mu} + \frac{z\delta}{\sqrt{n}} \right]$$

$$z\delta = \sqrt{\frac{1}{2} \ln \frac{2}{\delta}}$$

- $0 \leq \delta \leq 1$. The more confident we are the smaller the $\delta \rightarrow$ the larger the $Z_\delta \rightarrow$ the larger the confidence interval
 - The larger the n is, the smaller the confidence interval

Combining Arm Info

Combining Arm Info.

- The length of the bar, ε , is the confidence interval of μ_k .

How Many Samples?

How Many Samples?

Want to be $\epsilon/2$ accurate with probability $(1-\delta)/K$

$$\frac{\sqrt{\frac{1}{2} \ln\left(\frac{2K}{\delta}\right)}}{\sqrt{C}} \leq \varepsilon/2$$

$$\text{Quiz: } C \geq \frac{1}{\varepsilon} \ln \left(\frac{2k}{\delta} \right)$$

- C depends more on ε and less on δ (δ is in \ln)
 - PAC-style bounds for bandits.

Exploring Deterministic MDPs

MDP Optimization Criteria

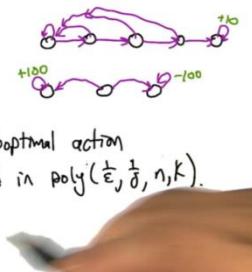
- ## 1. Explore randomly

7 Trap States

- ### 3 Mistake bound

The number of ϵ -suboptimal action choices is bounded in $\text{poly}(\frac{1}{\epsilon}, \frac{1}{\delta}, n, k)$.

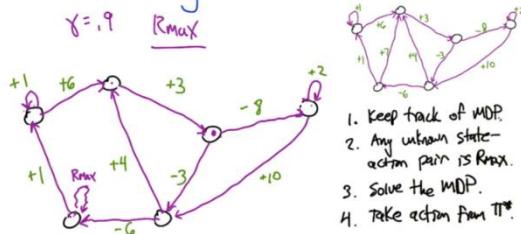
- #### 4. Minimize cheats



- In MDP, we can only choose action but not state
 - We can treat action selection as bandit problem (e.g. k action is a k arm bandit).

Exploring Deterministic MDPs

Exploring $\gamma = .9$ R_{max}



- **Rmax**: assume any state we have not fully explored (unknown state-action pair) has a self-loop of reward Rmax.
 - This way, the learner is encouraged to explore the whole MDP.

Rmax Analysis

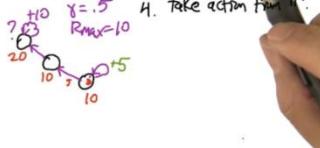
R_{max} Analysis

1. Once all edges known, no mistakes

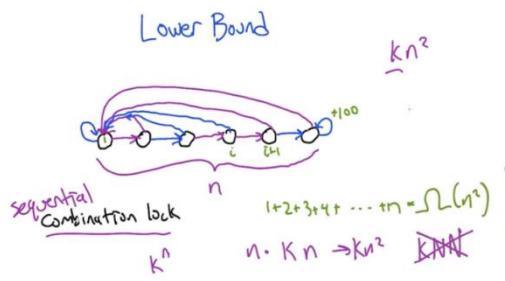
- ? Stop visiting unknown states?

2. $y = .5$ 4. Take action from "I".
? $\frac{1}{3}$ $R_{mg} \approx 10$

- 20 25



- The example: even with Rmax, there is a possibility that the learning stops visiting unknown states
 - Discount factor is important to determine if the agent is going to explore unknown edge.



- We have an upper bound saying the algorithm can solve the problem with n^2k steps
- now we are proving that no other algorithms can take less than n^2k steps
- The sequential combination lock has n states, in each state the learner can take k actions.
- at i th state, one action can send the learner to next state ($i + 1$), all other actions will take the learner to state 0, to get back to state i , the learner has to take at least i steps (the worst case deterministic MDP).
- In the worst case, it might take kn^2 steps to solve a deterministic MDP before stopping make mistakes.

General Rmax

General Stochastic MDPs

- stochastic (Hoeffding bound until estimates are accurate)
- sequential (unknown regions assumed optimistic)

For Stochastic properties, applying Hoeffding bound. For sequential deterministic MDP, Rmax can solve the problem. Now applying both, we can solve the stochastic MDP.

General Rmax

- Rmax
- Keep track of MDP. TR
 - Any unknown state-action pair is Rmax.
 - Solve the MDP.
 - Take action from π^* .
- $s, a \rightarrow s_{53} \frac{1}{c}$
- $\exists c: \text{Unknown } s, a$
if tried fewer than c times, then use ML estimate

- For Rmax algorithm, define Unknown s, a pair as: if tried fewer than C times. Then use maximum likelihood estimate.
- The C part is Hoeffding bound.

Some key ideas for this to work: Simulation Lemma and Explore-or-Exploit Lemma

Simulation Lemma

Simulation Lemma

Transitions and rewards off by α or less.

Fixed π ,

$$V_1^\pi(s) \quad V_2^\pi(s)$$

$$T_1, R_1 \quad T_2, R_2$$

$$\max_{s, a} |T_1(s, a, s') - T_2(s, a, s')| \leq \alpha \quad \max_{s, a} |R_1(s, a) - R_2(s, a)| \leq \alpha$$

- If we get an estimated MDP that is close to real MDP, then the optimizing the reward based on the estimate will be very close to the optimal policies for the real MDP.
- α is the difference between **estimated** and **real** MDP.
- α could be different if transition possibility and reward are not in the same scale, but we can always rescale things so that α is between 0 and 1.

Simulation Lemma

Define $G = \alpha + \gamma \alpha R_{\max}/(1-\gamma)$

$$E = |V_1(s) - V_2(s)| \leq \frac{G}{1-\gamma}$$

$$E = \frac{G}{1-\gamma} = \frac{\alpha + \gamma \alpha R_{\max}}{1-\gamma} = \frac{\alpha(1-\gamma)}{1-\gamma} + \frac{\gamma \alpha R_{\max}}{1-\gamma} = \frac{\alpha(1-\gamma)}{1-\gamma} + \frac{\gamma \alpha R_{\max}}{1-\gamma}$$

$$\frac{\alpha}{1-\gamma} = \frac{\alpha(1-\gamma)}{1-\gamma} + \frac{\gamma \alpha R_{\max}}{1-\gamma}$$

Hoeffding bound

- α is related to ϵ
- α is used to determine C
- Using **Hoeffding bound** and **union bound** to find a big enough C which can estimate an MDP which is close enough to real MDP.

Week 9: Generalization

Explore-or-Exploit Lemma

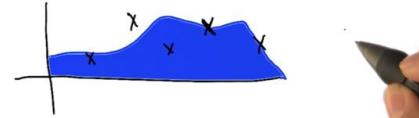
Explore or Exploit Lemma
 If all transitions are either accurately estimated or unknown, optimal policy is either near optimal or an unknown state is reached "quickly".
 $n \cdot K$ PAC

What Have We Learned?

What have we learned?

- bandits, decision making with randomness Hoeffding, union bound (learn how the world stochastic works)
- R_{\max} : optimism in the face of uncertainty. Explore at a distance. Grass is always greener. Sequential
- combined: R_{\max} + bandits stochastic + sequential
KWIK - know what it knows
- Hoeffding bound and Union bound let us know how certain we are about the environment
- R_{\max} makes sure we "visited" things enough to get accurate estimates

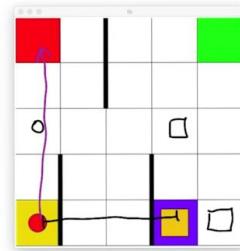
GENERALIZATION



Example : Taxi Representation!

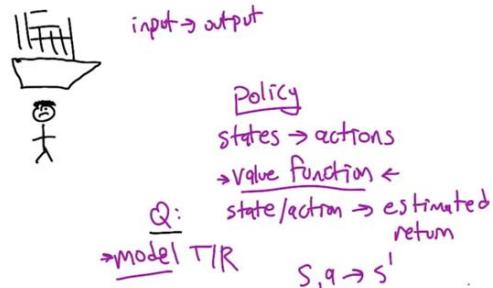
f_1 : low on fuel
 f_2 : x coordinate
 f_3 : y coordinate
 f_4 : near wall
 f_5 : always 1
 f_6 : ...

Inductive bias



- When there are a lot of states, we can use features to represent states (e.g. grouping states using some common features or using some base states to represent other states)
- Inductive bias: refer to the way algorithm prefer one solution against another.
- Representation can stay unchanged even when states changed (?)

Generalization Idea



- In RL, the goal is to learn some input-output mapping.
 - e.g. Policy maps state \rightarrow actions

- value function maps state/actions to expected return
- Model maps T/R (transition and reward)
- Generalization can happen on all these levels
- Most researchers are focused on **value function approximation**.

Basic Update Rule

General Form:

$$Q(s,a) = F(w^a, f(s))$$

$$\Delta w_i^a = \alpha (r + \gamma \max_{a'} Q(s', a') - Q(s, a)) \frac{\partial Q(s, a)}{\partial w_i^a}$$

$\underbrace{y^*}_{\text{Bellman equation}}$ $\underbrace{y}_{\text{TD error}}$ $\underbrace{\frac{\partial Q(s, a)}{\partial w_i^a}}_{\text{gradient step}}$
 bootstrapping TD error

- Q is represented by function approximator $F(w^a, f(s))$. w^a is some kind of weights, $f(s)$ takes a state s and generate a series of features.
- Experience tuple $\langle s, a, r, s' \rangle$ comes in, and do an update.
- TD error = the difference between predicted value (reward + discounted value of next state) and value of current state. TD error tells us the direction to move the parameters w . (the prediction is high, low or just right).
- α is learning rate.

Linear Value Function Approximation

Linear Value Function Approximation

$$Q(s,a) = \sum_{i=1}^n f_i(s) \cdot w_i^a \quad \leftarrow \text{parameters provide generalization}$$

Each action gets a set of n weights to represent the Q values for that action. Dot product.

Weights give "importance" of all the features in contributing to an action's value.

- Q value for an action is represented by features and a set of weights
- w^a_i are weights of features $f_i(s)$, which represent importance of features

Quiz

$$Q(s,a) = F(w^a, f(s)) = \sum_i w_i^a f_i(s)$$

$$\langle s, a, r, s' \rangle$$

$$\Delta w_i^a = \alpha (r + \gamma \max_{a'} Q(s', a') - Q(s, a)) \frac{\partial Q(s, a)}{\partial w_i^a}$$

$$\frac{\partial Q(s, a)}{\partial w_i^a} \Delta w_i^a = \boxed{f_i(s)}$$

Success and fail stories

Does it work?

Successes

3-layer backprop	: Tesar, Wang/Dietterich
CMAC	: Sutton <small>(computationally efficient)</small>
Linear nets	: Singh, Parr <small>(value-informative)</small>
Deep Nets	: Mnih et al. <small>Atari</small> <small>distilled to paper</small>
X	Tesar, Boyan/Moore → Sutton
Isbell	TD Gamma <small>need ad. work</small>
LSTT	Pollack - GA <small>if you do it right</small>
	<small>randomness</small> <small>→ MACE</small>

- 3-layer back prop: decision making in Jeopardy
- CMAC([wikipedia](#)) :

Baird's Counter Example

$\begin{matrix} & u_1 & u_2 & u_3 & u_4 & u_5 & u_6 & u_7 \\ 1 & 0 & 1 & 2 & 0 & 0 & 0 & 0 \\ 2 & 1 & 0 & 2 & 0 & 0 & 0 & 0 \\ 3 & 1 & 0 & 0 & 2 & 0 & 0 & 0 \\ 4 & 1 & 0 & 0 & 0 & 2 & 0 & 0 \\ 5 & 1 & 0 & 0 & 0 & 0 & 2 & 0 \\ 6 & 1 & 0 & 0 & 0 & 0 & 0 & 2 \\ 7 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{matrix}$ $1 \frac{1}{2} \frac{2}{3} \dots \frac{7}{8} \times c$	
<ul style="list-style-type: none"> - no rewards - deterministic - no choices - linear value function approximation - near tabular 	<ul style="list-style-type: none"> - value function representable - matrix representable - infinitely representable

Bad Update Sequence

$$V(1) = w_0 + 2w_1 = 3 \quad w_0 > 0 \quad w_1 = 1$$

$$V(2) = w_0 + 2w_2 \quad V(2) > V(1) \quad 8 > 3 \quad \checkmark$$

$$V(3) = w_0 + 2w_3 \quad \text{update rand robin} \quad \langle 1, 0, 7 \rangle$$

$$V(4) = w_0 + 2w_4 \quad \gamma = .9 \quad \langle 2, 0, 7 \rangle$$

$$V(5) = w_0 + 2w_5 \quad \alpha = .1 \quad \langle 3, 0, 7 \rangle$$

$$V(6) = w_0 + 2w_6 \quad \langle 4, 0, 7 \rangle$$

$$V(7) = 7w_0 + 2w_7 = 8 \quad \langle 5, 0, 7 \rangle$$

$$\langle 6, 0, 7 \rangle$$

$$\langle 7, 0, 7 \rangle$$

$$\underbrace{\Delta w_0}_{\langle 10, 7 \rangle} = \alpha (r + \gamma V(s') - V(s)) \frac{\partial V(s)}{\partial w_0}$$

Stealth quiz

$$= 0.42$$

- $V(s') = V(7) = 8$; $V(s) = V(1) = 3$, derivative is $w = 1$; so $\Delta w_0 = 0.1(0 + 0.9 * 8 - 3) * 1 = 0.42$;

- Δw_0 is 0.42, and $\Delta w_i, i = [0, 6]$ are all positive, so weight will increase with each update
- Δw_7 is negative when we update $\langle 7, 0, 7 \rangle$, but w_7 will still increase because its original value is $>> |\Delta w_7|$
- So, even the algorithm is doing the right thing, the updates will never converge

Bad Update Sequence

$$\begin{aligned} V(1) &= w_0 + 2w_1 \\ V(2) &= w_0 + 2w_2 \\ V(3) &= w_0 + 2w_3 \\ V(4) &= w_0 + 2w_4 \\ V(5) &= w_0 + 2w_5 \\ V(6) &= w_0 + 2w_6 \\ V(7) &= 7w_0 + w_7 \end{aligned}$$

$\Delta w_0 = \alpha \left(r + \gamma \sum_{s'} V(s') - V(s) \right) / \sum_{s'} w_s$

QUIZ:

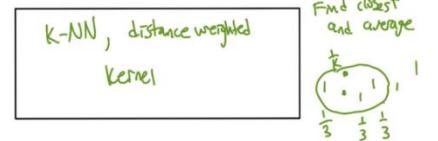
0	0	0	0	0	0	0	0
---	---	---	---	---	---	---	---

update round robin: $\langle 1, 0, 7 \rangle$
 $\gamma = .9$
 $\alpha = .1$

- Problem: given a certain points of value, estimate the function of line. How to estimates the points in-between?
- Naturally, weighted average of neighbors and get convex combinations of anchor points
- For points that are far from anchor points, are good estimate will be the average of the anchor points.
- More anchor points lead to less error in MDP estimation.

Averagers : QUIZ

Supervised learner
that can be expressed as an averager?



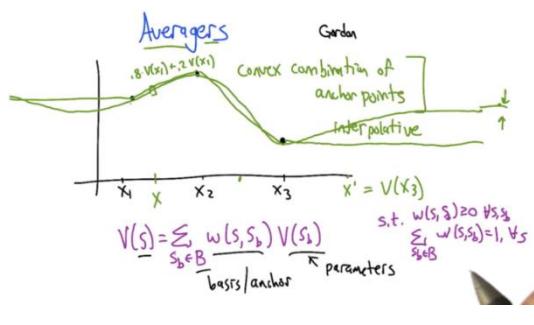
Connection to MDPs

$$\begin{aligned} V(s) &= \max_a \left(R(s, a) + \gamma \sum_{s'} T(s, a, s') V(s') \right) \quad \text{TD} \\ &= \max_a \left(R(s, a) + \gamma \sum_{s' \in B} \underbrace{\sum_{s_b \in B}}_{s' \in B} \underbrace{T(s, a, s_b)}_{s' \in B} V(s_b) \right) \\ &= \max_a \left(R(s, a) + \gamma \sum_{s' \in B} \left(\sum_{s_b \in B} T(s, a, s_b) \right) V(s_b) \right) \\ &\stackrel{\text{MDP!}}{=} \max_a \left(R(s, a) + \gamma \sum_{s' \in B} T'(s, a, s_b) V(s_b) \right) \end{aligned}$$

- In this case, the temporal differences are all zeroes because w_i is 0, which means no updates (changes) are going to happen.
Sticky
- If at any circumstance, when all values are zero, it could not escape the circumstance unless reward (r) is non-zero. **If reward is also zero, we can call the circumstance the solution of the MDP.**
- This holds when MDP is **deterministic**.
- The lesson: Shared weights can "not converge"!

- Using base states, we can represent the original MDP with MDP over S_b s

What have we learned?



- What have we learned?**
- Need for **ab** methods for generalization in RL
 - ↳ billions of states
 - ↳ supervised learning methods
 - linear function approximation
 - ① Value functions ② Policies ③ Models
 - gradient form
 - successes - TD gamma, Atari 2 new 2 no
 - problems - even linear need not work GTO 2 parallel converges linear loss functions
 - Averagers - convergence MDP unique value function Kestrel converges to MDP's unique Kestrel

- Need for generalization: some problems have billions of states; Method: applying supervised learning method.
- Methods: Linear function approximation.
 - value functions 2) policies; 3) models;
- Success: TD gamma, Atari, Fitted Q-iteration (but still have need not work cases)
- Problem cases: even linear need not work (baird counter problem)
- Averagers: can be viewed as MDP itself. Use anchor points to represent the whole MDP.
- LSPI: Least Squared policy iteration

- $O(S, Z)$ is the mapping function of Z to $\langle S, A \rangle$ (e.g. Given Z , what's the probability of S).

Quiz: POMDPs generalize MDPs

MDP: $\langle S, A, T, R \rangle$
 POMDP: $\langle S, A, Z, T, R, O \rangle$

$$Z = \begin{cases} S & \text{if } s = z \\ \square & \text{otherwise} \end{cases}$$

$$O(s, z) = \begin{cases} 1 & \text{if } s = z \\ 0 & \text{otherwise} \end{cases}$$

- When $O(S, Z) = 1$, the observable Z is S . So, Z and O together determine the relationship between S and Z in the POMDP.

Week 10: POMDPs, Partially observable MDPs

Partially Observable MDPs

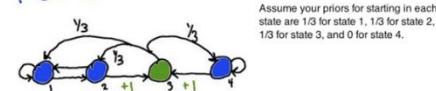
POMDP

POMDPs

A way of talking about "non-Markov" environments.
 MDP $\langle A, S, Z, T, R, O \rangle$ observables
 observation function
 $O(s, z)$ probability of seeing z in state s .

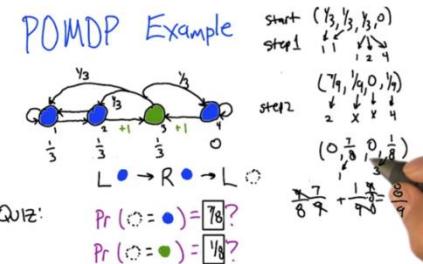
- POMDPs generalizes MDPs.
- In POMDP, MDP (represented by $\langle S, A, T, R \rangle$) is not directly observable to the agent. we can only observe the states. Z are observables to the agent, Z is a set of states in S .

POMDP Example



$L \bullet \rightarrow R \bullet \rightarrow L \circ$
 Quiz: $Pr(\circ = \bullet) = \frac{7}{8}?$
 $Pr(\circ = \bullet) = \frac{1}{8}?$

Solution



- Keep track of probability of ending up in each state after each step.
- There is a normalization after each step (not the right thing to do but works.).

State Estimation

(vector)

$$\begin{aligned}
 & b, a, z \rightarrow b' \quad \text{belief MDP} \leftarrow \text{POMDP} \\
 & \text{reward not observed} \\
 b'(s) &= \text{probability of being in state } s' \text{ after } b, a, z \\
 & = \sum_{s'} \Pr(s'|b, a, z) = \sum_{s'} \Pr(s|b, a, z) \cdot \Pr(z|b, a, z, s) \\
 & = \sum_s b(s) \Pr(z|s, a, z) \Pr(s|a, z) \\
 & = \frac{\sum_s b(s) O(s, z) T(s, a, z)}{\sum_{s'} b(s') O(s', z) T(s', a, z)} \quad \text{POMDP} \rightarrow \text{MDP (constant states)} \\
 & \text{LP, PI, VF} \quad \text{go}
 \end{aligned}$$

and α . α is a finite set of vectors. This is called piecewise linear and convex.

Piecewise Linear and Convex

Base case: There exists a set Γ_0 s.t.

$$V_0(b) = \max_{d \in \Gamma_0} d \cdot b = 0$$

QUIZ

$$\begin{aligned}
 \Gamma_0 &= \boxed{\{[0, \dots, 0]\}} \quad \vec{0} \\
 |\Gamma_0| &= 1
 \end{aligned}$$

$$\Gamma_0 = \{ [R(s_0, a), \dots, R(s_n, a)] \mid a \in A \}$$

- Γ_0 should be a vector of zeroes.

Piecewise Linear and Convex II

Assume we have Γ_{t-1} s.t. $V_{t-1}(b) = \max_{d \in \Gamma_{t-1}} d \cdot b$

Build Γ_t s.t. $V_t(b) = \max_{d \in \Gamma_t} d \cdot b$

$$V_t(b) = \max_{a \in A} V_t^a(b)$$

$$V_t^a(b) = \sum_z V_t^{a,z}(b)$$

$$V_t^{a,z}(b) = \sum_s R(s, a) \underbrace{b(s)}_{|z|} + \gamma \Pr(z|b, a) V_t^{(SE(b, a, z))}$$

$$\begin{aligned}
 b'(s') &= \Pr(s'|b, a, z) \\
 &= \Pr(z|b, a, s') \Pr(s'|b, a) / \Pr(z|b, a) \\
 &= \Pr(z|b, a, s') \sum \Pr(s'|s, b, a) \Pr(s|b, a) / \Pr(z|b, a) \\
 &= O(s', z) \sum_s T(s, a, s') b(s) / \Pr(z|b, a)
 \end{aligned}$$

- Note: belief MDP has infinite number of belief states which make VI, LP, PI impossible because they can only deal with finite number of states.

Value Iteration in POMDPs

Value Iteration in POMDPs

$$V_t(b) = \max_a (R(b, a) + \gamma \sum_z \Pr(z|b, a) \cdot V_{t-1}(b'))$$

where $b' = SE(b, a, z)$

CNIM: $V_t(b) = \max_{d \in \Gamma_t} d \cdot b = \max_{d \in \Gamma_t} \sum_s b(s) \cdot d(s)$

Sondik

- Piecewise linear means $V_t(b)$ can be represented by $V^a_t(b)$, $V^a_t(b)$ can be represented by $V^{a,z}_t(b)$ and $V^{a,z}_t(b)$ can be represented by something like a Bellman equation

Piecewise Linear and Convex III

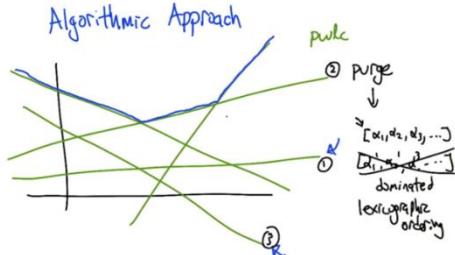
$$\begin{aligned}
 \Gamma_t &= \bigcup_a \Gamma_t^a \\
 \Gamma_t^a &= \bigoplus_z \Gamma_t^{a,z} \\
 \Gamma_t^{a,z} &= \{ \sum_s R(s, a) + \gamma \sum_{s'} \alpha(s') O(s', z) T(s, a, s') \mid \alpha \in \Gamma_{t-1} \} \\
 V_t^{(SE(b, a, z))} &= \max_{\alpha \in \Gamma_{t-1}} \sum_s \alpha(s) O(s, z) \sum_{s'} T(s, a, s') b(s') \\
 |\Gamma_t| &= |\Gamma_{t-1}|^{|z|} \cdot |A| \quad \rightarrow \text{pomdp-solve}
 \end{aligned}$$

- In a similar sense, Γ_t can be represented by Γ_t^a , $\Gamma_t^{a,z}$ and something like a Bellman equation.
- Thus the belief MDP becomes computable.
- Pomdp-solve is an existing solver for pomdp.
- Guaranteed to be exponential because size of Γ is exponential in respect to the size of observables: $|\Gamma_t| = |\Gamma_{t-1}|^{|z|} \cdot |A|$

- Initialized $V_0(b)$ as 0
- $V_t(b)$ can be written as something similar to bellman equation
- Claim $V_t(b)$ can be represented by the maximum of a set of linear functions of b

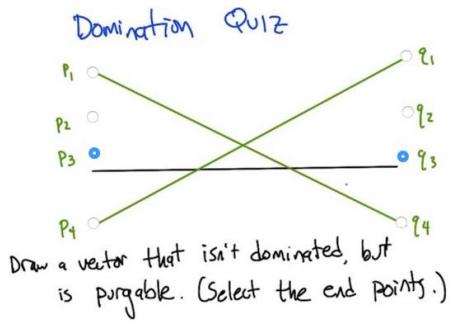
- The calculation could be faster if Q function ($V^a_t(b)$ is the Q function) has a small representation.

Algorithm



In the figure:

- All the linear functions (vectors) are dominated by the blue lines which involves three linear functions.
- Function ① and ③ can be purged because they have no contribution to the overall vector



- Domination: if one vector is always less than (in this case) another vector, we say it's dominated by that another vector.
- Vector is purgeable if it does not contribute to the max

Learning POMDP

RL for POMDPs

Planning - Know everything (Model)
RL - don't know, need to interact

exact optimal answer
undecidable
↳ life is profound

- Model-based RL \leftarrow learned model & use it
- Model-free RL \leftarrow don't bother

learn POMDP
map observations to actions

- Planning: if know the model of POMDP, can run value iteration. Since POMDP have infinite number of states, we can only get approximation of optimal policy.
 - VI may not be converge
 - PI is not possible because infinite believe states
- Difference between planning and reinforcement learning:
 - Planning: know the model - planning in POMDP is undecidable.
 - RL: we don't know the model have to interact to learn
- RL has 2 types:
 - Model based RL: learn the model (POMDP) and use it
 - Model free RL: don't learn the model but mapping observables and action.

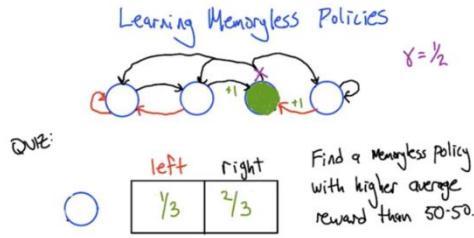
Learning a POMDP

Q1T	uncertain		controlled	
	observed	partially observed	MC	MDP
	HMM	POMDP	mixture of Gaussians HMM POMDP chromosome	
E			expected value of hidden variables	M

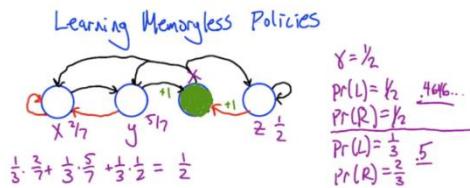
- MC: Markov Chain, observed but uncontrolled; HMM: Hidden Markov Model, partially observed and uncontrolled; MDP: observed and controlled. POMDP: partially observed and controlled.

- EM: *Expectation Maximization*, an procedure to learn the hidden property in HMM or POMDP.
- Model based Learning.

Learning Memoryless Policies (Model free RL of POMDP)



- In the MDP, There 3 blue states, to get to the green state, the agent needs to move to right in 2 of the states and move to left in another state. So the solution for this quiz is $1/3$ and $2/3$.



- Solve and compare the value function based on the 50-50 probability (policy) and $1/3$ - $2/3$ policy.

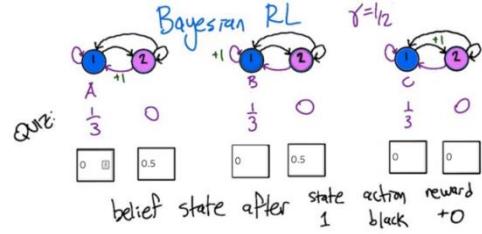
$$X = p * (0.5 * X) + (1 - p)(0.5 * Y) \Rightarrow X = 0.5 * (1 - p)Y / (1 - 0.5p)$$

$$Y = p * (0.5 * X) + (1 - p) \Rightarrow X = 2 * (Y - 1 + p) / p$$

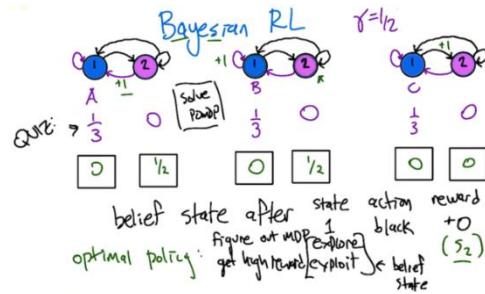
$$Z = p * (1) + (1 - p) (0.5 * Z) \Rightarrow$$

$$V = 1/3 * (X + Y + Z)$$

RL as POMDP:



- We can be in MDP A or B but not C, and we will be in state 2.
- Next, we can take a purple action to figure out which MDP we are in.



- So the optimal policy is figuring out which MDP we are in and get hight reward by exploiting in the belief state.
- RL itself is like POMDP which requires exploring and exploiting.

Bayesian RL
 too expensive
 RL is actually planning in a kind of hidden state parameters of MDP
 piecewise Polynomial and convex
 BETTER
 LEarning Everything
 Poupart et al.

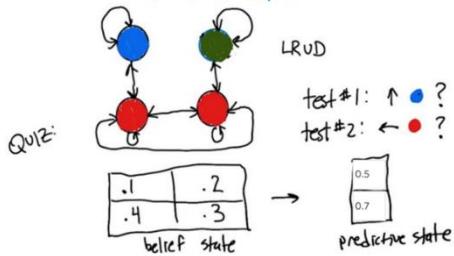
Predictive State Representation

Predictive State Representation

POMDP: belief state is distribution over states.
states never observed. learning
Do they even exist?

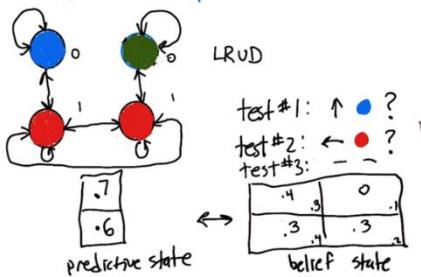
PSR: predictive state is probabilities of future predictions.

PSR Example



- The POMDP has 4 states but 3 observables represented by colors where the two red states are not distinguishable.
- We can do test to get predictions of states. (e.g., a test can be moving up, and see the probability of ending up at the blue test).

PSR Example



- Belief state and predictive state does not have a one-on-one mapping.
- Increasing number of tests will constrain the mapping.

PSR Theorem

Any n-state POMDP can be represented by a PSR with no more than n tests, each of which

is no longer than n steps long.
test #1: ↑ o ↑ o ?
test #2: ↑ o ?
Holmes & Isbell



easy

hard

continuous

single

Why go to PSR?

- We cannot observe all the states in POMDP and they might not even exist.
- We can do some test to figure out what state we might be in.
- Learning a PSR sometimes is easier than learning a POMDP.
- Number of test is always going to be no more than the number of states

What Have We Learned?

POMDP states = predictors

↳ belief states
↳ generalize MDPs
↳ hard (to accept)
↳ Value iteration - pwlc functions (sets of linear functions)
RL hard
↳ Expectation Maximization (EM)
↳ Memoryless may be random
Bayesian RL Planning/learning frequentist RL
Predictive state representation (PSR) exploitation is exploration "incentive problem" (sutton)

Week 11: Generalizing Generalization

GENERALIZING

GENERALIZING

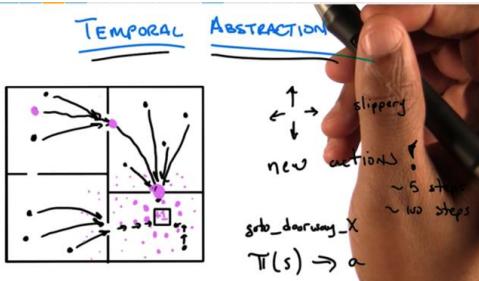
abstraction
scale
search

WHAT MAKES RL HARD?

- delayed reward \rightarrow not like SL
- bootstrapping \rightarrow need exploration
- # states ; # actions
 \rightarrow function approx
 $\& V(s), Q(s,a)$

- Delayed reward: agent has weak feedback, the reward is a moving target
- Need exploration to learn the model or action-reward pair for all or a good number of states.
- Computationally, the complexity of RL depends on number of states and # of Actions.
 - Function approximation over value function ($V(s)$, $Q(s,a)$) is abstraction over states, not actions.
 - The focus of this class is abstraction over actions

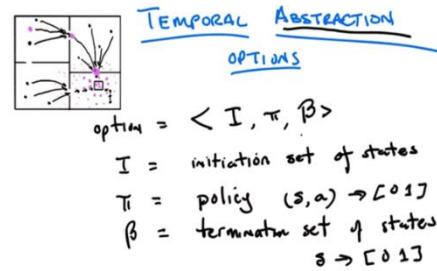
Temporal Abstraction



- Original grid world problem: walls separated big square into four rooms, agent could be in any location. Agent's actions are move L,R,U,D and its goal is to reach the small square in the bottom left room.
- A set of new actions (goto_doorway_X, X is orientation) can be generated to represent original actions.
- **Temporal Abstraction** is representing many actions with one or a few actions.(without doing anything to the states)

- Temporal Abstraction align many actions into one and thus makes a lot of states equivalent.
- Temporal Abstraction has computational benefits.

Temporal Abstraction Options



Option is a tuple $\langle I, \pi, \beta \rangle$

- I is **initiation** set of states
- π here is policy: probability of take action a in states s , $(s, a) \rightarrow [0 1]$
- β is **termination** set of states and it's a set of probability of ending at state s , $[0 1]$.

Temporal Abstraction Option Function

$$V_{t+1}(s) = \max_o (R(s,o) + \sum_{s'} F(s,o,s') V_t(s'))$$

$$R = E[r_1 + r_2 + \dots + r_k | s, k \text{ steps}]$$

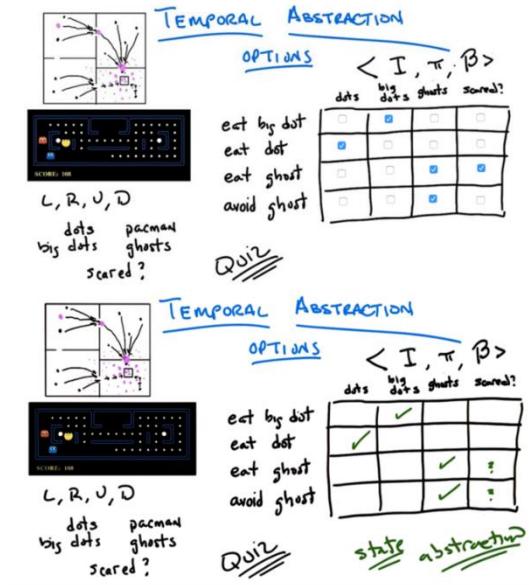
$$F = E[\gamma^k \Delta_{s_k, s'} | s, k \text{ steps}]$$

MDP: $\xrightarrow{\text{atomic actions}}$...
SMDP: $\xrightarrow{\text{variable time actions}}$

- The function is a rewrite of Bellman function
 - Using "o" to replace "a" (O is the generalization of A).
 - $V(s)$ (and $V(s')$) is the value function that needs to update, $R(s,o)$ is reward in s and choose a , F works like a transition function. The discount factor is hidden.
- Using options kind of violates the temporal assumption of MDPs

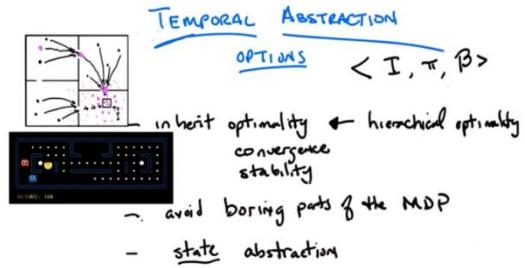
- MDPs have atomic actions, reward can be easily discounted for each step.
- Using options ends up with variable time actions, discount factor is hidden.
- If o represent k steps, R and F are actually discounted. this is Semi-MDP or SMDP
 - We can turn non-Markovian stuff into Markovian by keep tracking of its history.

Pac-Man Problems



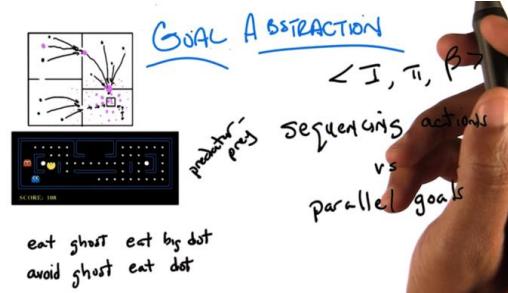
- We can learn two things from the example:
 - If done improperly, temporal Abstraction might not reach optimal policy.
 - Temporal Abstraction might introduce state abstractions (reduce the state space) so the problem can be solved more efficiently.

How It Comes Together

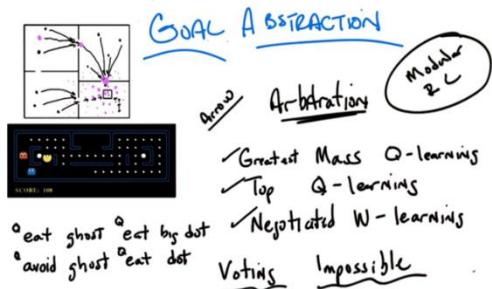


- We can see options and state representations as high level representation. In fact, actions and states are also somewhat made up. Agent's goal is to make decisions with respect to those descriptions of the world, no matter they are action or option, states or abstract states
- If construct options smartly, we might be able to ignore some states (decrease the state space) to reduce computational resource requirement.

Goal Abstraction



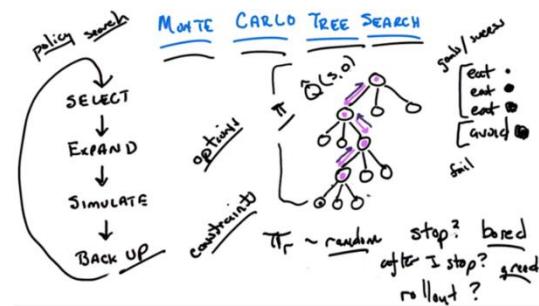
- Agent has multiple parallel goals (predator-prey scenario), at any given time, only one goal is dominating (more important).
- β (the probability of terminated in a state) can be seen the probability of accomplishing a goal or the probability of another goal becomes more important (switching goals).
- Options give ways to think actions as accomplishing parallel goals.
- The goals do not have to be in the same scale.



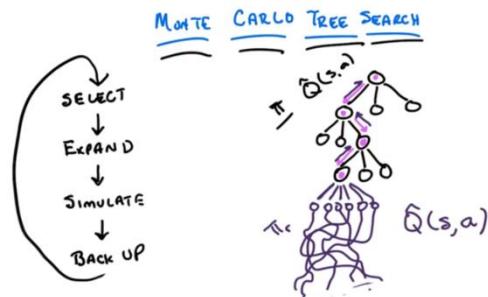
- Modular RL sees options as sub-agents with different goals. There are three ways of choosing goals:
 - Greatest mass Q-learning: each action as a Q function. For each option, we sum up the Q functions of each action in it. The option with largest Q get executed. (Might end up with agent with small Q on every action).
 - Top Q-learnings: choose the action with highest Q. (but the agent might put high Q on many actions)
 - Negotiated W-learning: Minimize loss.
- Modular RL is often impossible because a fair voting system is hard to construct. e.g. the modules might have incompatible goals.

the policy at unknown state. Then repeat the "Select, Expand, simulate, back up" process.

- Initial policy π , we can **SELECT** actions following it. it can be updated at each iteration of tree search.
- Rollout policy π_r , we can **simulate** actions and sample from them starting from the unknown state following it.
- We **back up** after stimulation with the simulated result to update π .
- We figured out what action to take at the unknown state and states above and **expand** π to the previous unknown state.
- Now we can repeat the process to search deeper.



Monte Carlo Tree Search



In the figure above, circles are states, edges are transitions. $\pi = Q^*(s,a)$ is the policy of the known part of the tree. In these states, we know what action to take following π (pink edges). When reach an unknown state, we apply the rollout policy π_r , and simulate actions to take deep in the tree, and then we backup and update π_r and π to figure out what to select at each state, including the unknown state where we started the simulation. π gets expanded as we figure out

- When to stop? Learn deeper enough before reach the computational resource cap
- Rollup policy π_r can be random. We know the action is good because I get better result by behaving randomly from that point on.
- Instead of purely random, one can behave randomly in respect to constraints. (E.g. not eaten by ghost).
 - constraints: defined by failure
 - goals: defined by success
- MCTS is compatible with options to perform the tree search. in this case, $\pi = Q_{\text{hat}}(s,o)$;
- The Monte Carlo Tree Search can be seen as policy search. When reach a state we are

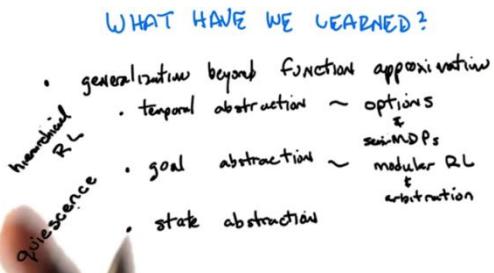
not confident, an inner loop is executed to do some RL.

MCTS Properties

- MONTE CARLO TREE SEARCH
- useful for large state spaces
- need lots of samples to get a good estimate
- planning time independent of $|S|$ but depends on X
- running time exponential in the horizon $\mathcal{O}(|A|^X)^H$

- MCTS is useful for large state spaces and need lots of samples to get good estimates
- Planning time is independent of $|S|$
- Running time is $\mathcal{O}(|A|^X)^H$.
 - X is how many steps does the search need to go through
 - $|A|$ is the size of the action space
- The tradeoff is between number of states and how far away are we going to search

Recap

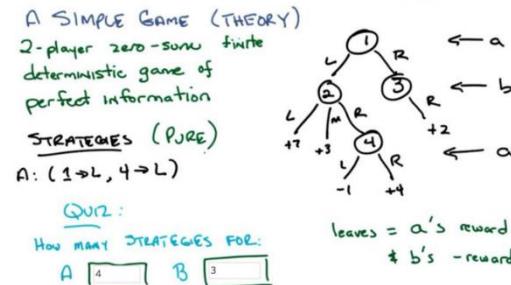


Week 12: Game Theory I

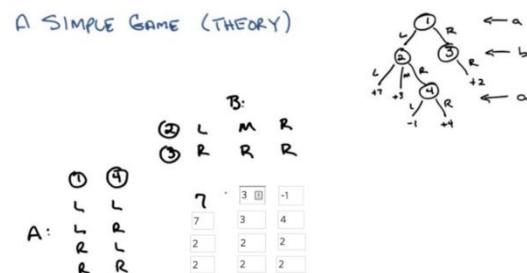
- GAME THEORY
- MATHEMATICS OF CONFLICT
 - SINGLE AGENTS → MULTIPLE AGENTS
 - ECONOMICS (+ POLITICS ...)
 - INCREASINGLY A PART OF AI/ML

- Game theory is mathematics of conflict of interests.
- It generalizes the RL from single agent to multiple agents.
- It is of the interest of economics, politics, and sociology or even biology since these fields often deal with agents with many agents with conflicts of interests.

Example



- The example game is represented as a tree, nodes are states, edges are transitions and leafs are rewards.
- The example game has 2 agent (a and b), all the rewards added up to be a constant (zero-sum), no stochastic transitions. The leafs show the reward that agent a can get and b gets (-1 * reward).
- Strategies: what action should the agent take at each state it could be in. (equivalent to policy in RL)
- The agents are assumed to be rational.

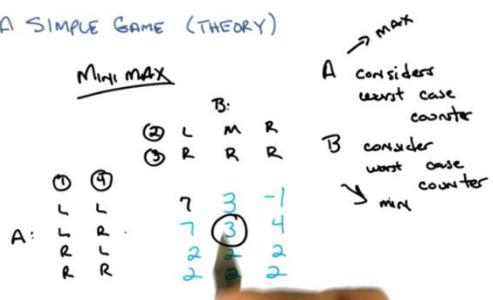


- This is a simple example of 2-agent zero-sum game

- Zero-sum means the reward of A and B will sum to 0 for any strategy.
- A matrix is enough to represent a game.
- Once the matrix is here, the tree doesn't matter now.
- Need to learn to generate the matrix from the tree.

Minimax

A SIMPLE GAME (THEORY)



- A and B has the same goal, maximize their own reward (and minimize others reward).
- If A and B behave rationally, they will reach the same strategy.

FUNDAMENTAL RESULT

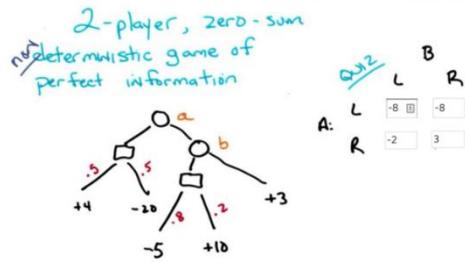
In a 2-player, zero-sum deterministic game of perfect information

Minimax ≡ Maximin... and there always exists an optimal pure strategy for each player

This is important so I am writing it down:

- In a 2-agent, zero-sum, deterministic game of perfect information, Minimax ≡ Maximin,
- And there always exist an **optimal pure strategy** for each agent.
- Based on the strategy matrix, one can always build at least one tree.

Now, to make the problem a bit more complex, we change the game to be non-deterministic:



- Introduce the **chance box**, so that transition is non-deterministic.
- Construct the strategy matrix based on tree (but could not reconstruct tree based on matrix)
 - Note: from the matrix, we don't know if the tree is deterministic or not.
- The minimax theorem (Von Neumann theorem) still holds
 - Minimax ≡ Maximin
 - Optimal pure strategy exists.

Minipoker

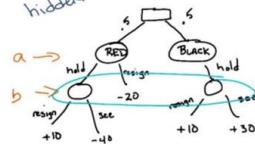
2-player, zero-sum non-deterministic game of perfect information



- A is dealt a card, red or black 50%
- A may resign if red: -20 cents for A else A holds
- B resigns: +10 cents
- B sees:
 - if red: -40 cents
 - if black: +30 cents

2-player, zero-sum non-deterministic game of perfect information

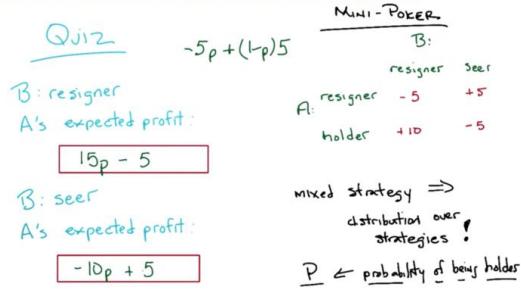
QUIZ



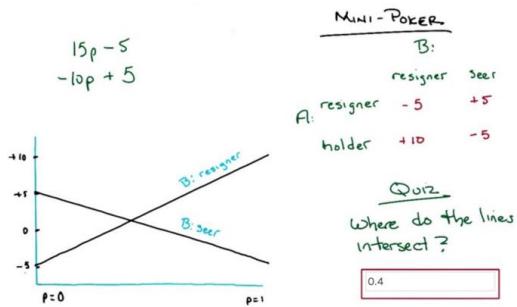
- In the minipoker game, B will not know all the information, so it's a 2-agent, zero-sum, non-deterministic game of **hidden information**

- In this case, Minimax \neq Maximin and there will be no optimal pure strategy.
- But there will be mixed Strategy, which is a distribution of strategies.

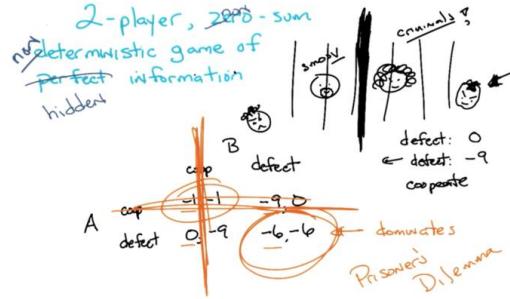
Mixed Strategy



- A's expected strategy are linear equation, which can be represented by lines.



- The mixed strategy should be at the intercept of the two lines in this case.
- If the two lines both have positive slope, the mixed strategy should be at the far right; if negative slope, the strategy should be at the far left.
- The expected value of the game is deterministic still.
- Although the strategy is mixed, there is still an expected value, in this case, the expected value is when p is 0.4, and value is 1.



- Now, we are making the game non-zero-sum.
- The prisoner's dilemma is a 2-agent, non-zero-sum, non-deterministic game of hidden information
- Assume the agents are rational, both of them should defect.

A Beautiful Equilibrium

A BEAUTIFUL EQUILIBRIUM
n players with strategies s_1, s_2, \dots, s_n
 $s_i^* \in S_1, s_2^* \in S_2, \dots, s_n^* \in S_n$
are a NASH EQUILIBRIUM iff
 $\forall i, s_i^* = \underset{s_i}{\operatorname{argmax}} \text{utility}_i(s_1^*, \dots, s_{i-1}^*, s_i^*, s_{i+1}^*, \dots, s_n^*)$
no reason for any one to change

- In practice, if we are in a Nash Equilibrium, any agent will have no good reason to change strategy (pure or mixed).

A BEAUTIFUL EQUILIBRIUM

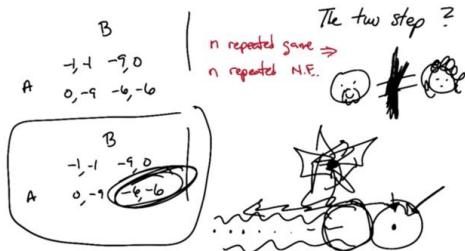
- In the n-player pure strategy game, if elimination of strictly dominated strategies eliminates all but one combination, that combination is the unique N.E.
- Any N.E. will survive elimination of strictly dominated strategies
- If n is finite & $\forall i, S_i$ is finite \exists (mixed) N.E.

- In the n-player pure strategy game, if elimination of strictly dominated strategies eliminates all but one combination, that combination is the unique NE.

Snitch

- Any N.E. will survive elimination of strictly dominated strategies
- If n is finite, for each set of finite strategies, then there will be at least one strategy is N.E.

A BEAUTIFUL EQUILIBRIUM ?



- What if playing the game multiple times?
- Only the last game matters \rightarrow the last game will be N.E. \rightarrow since the last game is known now, the last game moves to the game before it. \rightarrow the last game will be N.E. \rightarrow repeat.... \rightarrow all the game should will be N.E.

Recap

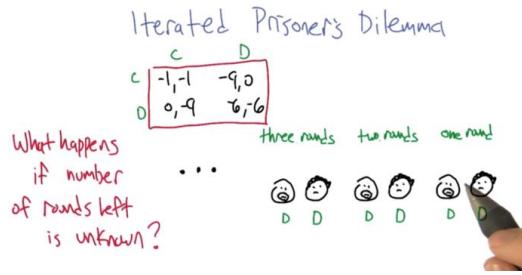
GAME THEORY: WHAT HAVE WE LEARNED?

- STRATEGIES: PURE vs MIXED
- ANDREW MOORE
- PRISONER'S DILEMMA
- NASH !
- MECHANISM DESIGN
- WE ARE ALL IN THE MATRIX
- MINIMAX
- HIDDEN / PERFECT
- (NON) ZERO SUM
- (NON) DETERMINISTIC

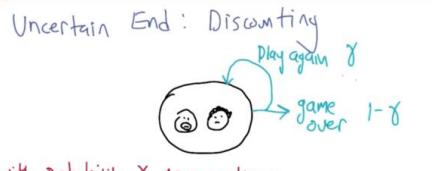
[Andrew Moore's slides on Zero-Sum Games](#)
[Andrew Moore's slides on Non-Zero-Sum Games](#)

Week 13: Game Theory II & III

11B Game Theory ReLoaded



- Prisoner should always defects at the last game (say we have k games), then the k th game is known
 - If you do the above choice, then the last game is known, and the last unknown game becomes to be the one before.
 - Since prisoner should defect in the last unknown game, then it should defect on the $(k - 1)$ th game. And the $(k - 1)$ th game becomes known.
- So prisoners should always defect at every each game
- But what happens when the last a few games are unknown?



With probability γ , game continues.
 Every round could be your last. or not.
 Expected # rounds? $\gamma = .99 \rightarrow$ rounds 100 $\frac{1}{1-\gamma}$
 (Finite if $\gamma \leq 1$)

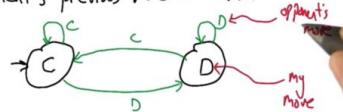
- the probability that the game will continue is γ
- Then the expected number of rounds is $1/(1 - \gamma)$.

Tit For Tat

Famous IPD strategy

→ cooperate on first round

→ copy opponent's previous move thereafter



- Finite representation of an unbounded number of rounds
- Copy opponent's previous move thereafter.

Tit For Tat

QUIZ: What does TFT do when playing against these strategies?

strategies	always defect	always cooperate	C-D-D:D... C-D-C-D...
always defect	□	□	□
always cooperate	□	□	□
TFT	□	□	□
D-C-D-C-D...	□	□	□

What to do if we Facing TFT?

Facing TFT

What's the best response to TFT?

always D $0 + \frac{-6}{1-\gamma}$ best for low γ

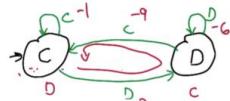
always C $\frac{-1}{1-\gamma}$ best for high γ

For what γ are they equally good? $\boxed{0.1667}$

C	D
C	-1, -1 -9, 0
D	0, 9 6, -6

Best Response To A

Finite-state Strategy



States labeled with opponents choice

edges labeled with our choice

edges annotated with our payoff for that choice

our choice impacts our payoff and
future decisions of the opponent $\textcircled{1}$ MDP! $\textcircled{2}$ Value Iteration
always C, always D, D-C-D-C-D-... =

- TFT is like an MDP, we can solve TFT by solving MDP
- If our opponent using TFT, then our decision will affect the opponent's decision.

Best Responses in IPD

What's the best response to each strategy? ($\gamma > \frac{1}{6}$)

opponent strategy	always C	always D	TFT	Possible response
always C	□	□	□	
always D	□	□	□	
TFT	□	□	□	

Best Responses in IPD

What's the best response to each strategy? ($\gamma > \frac{1}{6}$)

opponent strategy	always C	always D	TFT	Possible response
always C	✓	□	✓	
always D	✓	□	✓	
TFT	✓	✓	✓	cooperative!

Mutual best response. Pair of strategies where each best response to other. Nash!

- Playing the game multiple times (unknown times) opens up the possibility of another Nash Equilibrium which is cooperating with each other.
- In this case, both agents playing TFT is a best way to respond to each other.
- But what if γ is not $1/6$ here??? Or is γ matter?

Folk Theorem

Repeated Games and the Folk Theorem

General idea: In repeated games, the possibility of retaliation opens the door for cooperation.

What's a "Folk Theorem"?

In mathematics: Results known, at least to experts in the field, and considered to have established status, but not published in complete form.



Folk Theorem

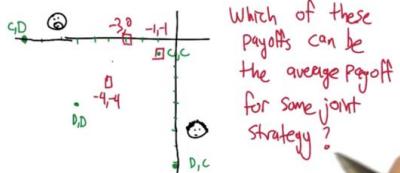
In game theory though, Folk Theorem refers to a particular result:

Describes the set of payoffs that can result from Nash strategies in repeated games.

- Describes the set of payoffs that can result from Nash strategies in repeated games.

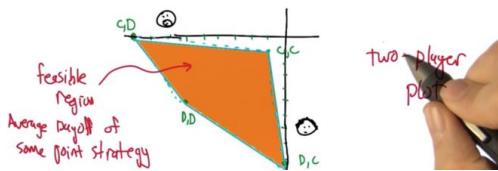
Repeated Games and the Folk Theorem

General idea: In repeated games, the possibility of retaliation opens the door for cooperation.



Repeated Games and the Folk Theorem

General idea: In repeated games, the possibility of retaliation opens the door for cooperation.



- MinMax Point: DD point
 - Feasible payoff: yellow filled region
 - Acceptable payoff: any point that to the right-above of MinMax point
 - Feasible acceptable payoff: intersection of 1 and 2.

Folk Theorem

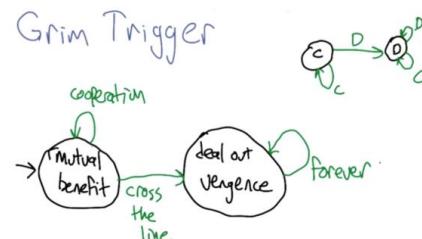
Folk Theorem

Any feasible payoff profile that strictly dominates the minmax/security level profile can be realized as a Nash equilibrium payoff profile, with sufficiently large discount factor

Proof: If it strictly dominates the min max profile, can use it as a threat. Better off doing what you are told!

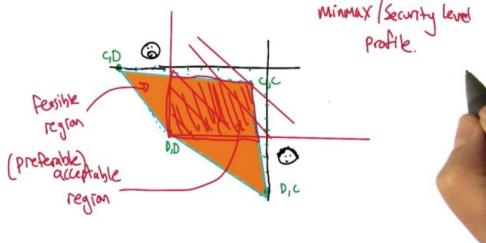
- Any feasible payoff profile that strictly dominates the minmax/security level profile can be realized as a Nash equilibrium payoff profile with sufficiently large discount factor.
 - Translation: do what you are told (feasible payoff), or you will be punished (minMax payoff profile). See the next slide for graph representation.

Proof of Folk Theorem



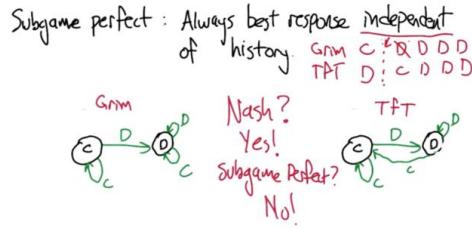
- x is the probability that Smooth chooses "B", Then $(1 - x)$ is the probability that Smooth chooses "S". So Curly will get x if he chooses "B", and $2(1 - x)$ if he chooses "S". Solving $x = 2(1 - x)$ will get x .

Repeated Games and the Folk Theorem

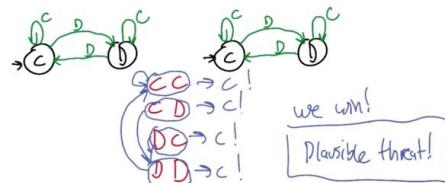


- black color--> actions of one agent
 - green color--> another agent
 - Grim Trigger: once cross line, will always deal with default.

Implausible Threats



Pavlov is Subgame Perfect!



Average reward: mutual cooperation

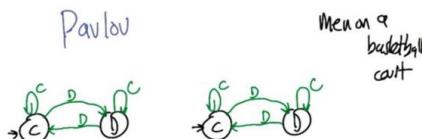
- Subgame perfect equilibrium: means taking test response independent of history.
- If under a threat the machine can reach a Nash equilibrium but cannot reach subgame perfect, then it is an implausible threat (the agent won't carry out the threat)

Is TFT vs. TFT subgame perfect?

- Yes
 No Here is a sequence that proves it:



- If #1 is defecting, #2 will get better expected reward if it keeps cooperating instead copying the action of one. SO TfT vs. TfT is NOT subgame perfect.



TFT: repeat what opponent did
 Pavlov: cooperate if agree, defect if disagree
 Is Pavlov v. Pavlov Nash? Yes No

- Pavlov can reach Nash equilibrium if they both choose to cooperate.

- Pavlov is subgame perfect and reaches NE. So the threat is plausible!

Computational Folk theorem
 2 player bimatrix game
 \rightarrow average reward

Can build Pavlov-like machines for any game.
 Construct subgame perfect Nash equilibrium
 for any game in polynomial time.

- Pavlov if possible
- zero-sum like (Solve an LP)
- at most one player improves

Peter Stone
 Me 3
 MIMIC

Stochastic Games and Multiagent RL

Stochastic Games and Multiagent RL

- Nash?
 pair of policies
 st. neither would
 prefer to switch.
 A: $y_2 \rightarrow y_2$
 B: $1 \rightarrow y_2$
-

MDP: RL :: stochastic game: multiagent RL

- Game description: A and B are two agents. If they go north, there is a semi-wall with 50% chance of going through. If they both go to the center, they will flip a coin to decide who will take the block. When one agent reaches \$, game ends and the agent who reaches \$ gets 100 reward.
- NE exist? Yes, both go center.

Stochastic Games (Shapley)

S : states s

A_i : Actions for player i . $a_1, b_1 \in A_1, b_2 \in A_2$

T : Transitions $T(s, (a, b), s')$

R_i : Rewards for player i $R_i(s, (a, b))$, $R_2(s, (a, b))$

γ : Discount

- Stochastic Games is generalization of both MDP and single agent RL.

Models & Stochastic Games

$\langle S, A_i, T, R_i, \gamma \rangle$

- ① $R_1 = -R_2$
 - ② $T(s, (a, b), s') = T(s, (a, b'), s) \quad \forall b'$
 $R_2(s, (a, b)) = 0, R_1(s, (a, b)) = R(s, (a, b')) \quad \forall b'$
 - ③ $|S| = 1$
- (A) MDP, (B) zero sum stochastic game
(C) repeated game

- If put some constraints on the stochastic games, we will get some of the games/mdp etc we learned before.
- Since stochastic games are generalized MDP, so we can try to generalize the method to solved the MDP to solve stochastic games.

Zero-Sum Stochastic Games

Zero-Sum Stochastic Games

$$Q_i^*(s, (a, b)) = R_i(s, (a, b)) + \gamma \sum_{s'} T(s, (a, b), s') \max_{a'_1, b'_1} Q_i^*(s', (a', b'))$$

$$\langle s, (a, b), (r, f), s' \rangle: Q_i(s, (a, b)) \leftarrow r_i + \gamma \max_{a'_1, b'_1} Q_i(s', (a', b')) \quad \text{Nash-Q}$$

- value iteration works
- minimax-Q converges
- unique solution to Q^*
- policies can be computed independently
- update efficient
- Q functions sufficient to specify policy

- Minimax can be used in a Bellman-like equation.

- VI works for the situation, minimax-Q converges, and Q^* is unique, policies can be computed.

General-Sum Games

General Zero-sum Stochastic Games

$$Q_i^*(s, (a, b)) = R_i(s, (a, b)) + \gamma \sum_{s'} T(s, (a, b), s') \max_{a'_1, b'_1} Q_i^*(s', (a', b')) \quad \text{Nash}$$

$$\langle s, (a, b), (r, f), s' \rangle: Q_i(s, (a, b)) \leftarrow r_i + \gamma \max_{a'_1, b'_1} Q_i(s', (a', b')) \quad \text{Nash} \quad \text{Nash} \quad \text{Nash} \rightarrow Q$$

- value iteration works doesn't work

- minimax-Q converges doesn't converge

- No unique solution to Q^*

- Policies can't be computed independently incompatible

- update not efficient $P = PBM$ insufficient

- Q functions not sufficient to specify policy

- When we general-sum case, Q type of algorithms will not work.
- The problem is not solved but there are a lot of ideas about =

Lots of Ideas

- repeated stochastic games (folk theorem)
- cheap talk → correlated equilibria
- cognitive hierarchy → best responses
- side payments (cooperative values)

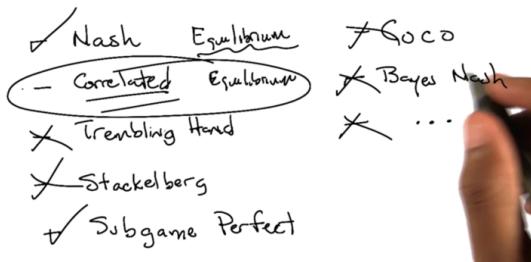
Recap

What Have we Learned?

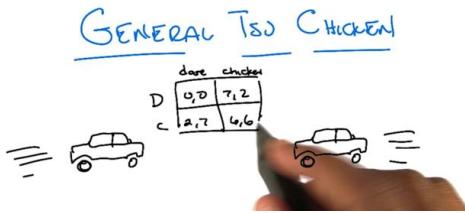
- Iterated PD
- connect IPD & RL (discounting) repeated games
- folk theorem (threats)
- subgame perfection, plausible threats
- computational folk theorem more acceptable
- stochastic games, generalize MDP, repeated games
- zero sum stochastic games. minimax Q works.
- general sum games. Nash Q doesn't. (End hopefully)

11C Game Theory Revolutions

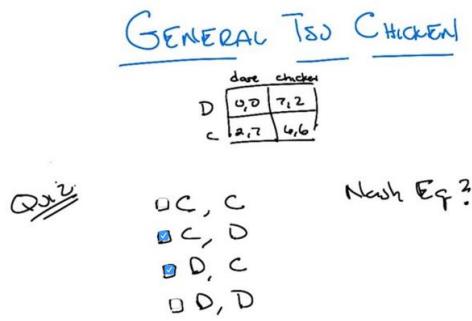
SOLUTION CONCEPTS



Correlated equilibrium

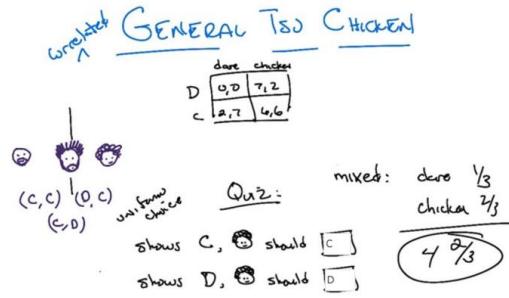


- In the chicken game, if both *Dare*, both get 0 reward (dies), if both *chicken*, both get 6 rewards. If one dare and the other chicken out, one get 7 and the other get 2 (still alive).

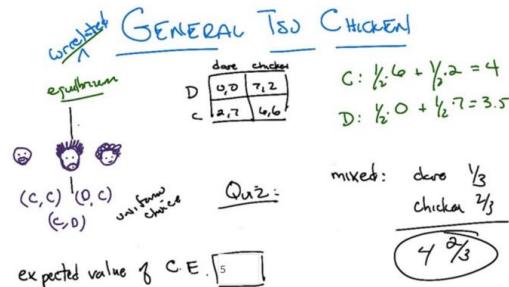


- CC is not NE because if I know you will be C, I will change to D to get more reward.
- CD or DC is NE because C is not willing to change to D, and D is not willing to change to C to get lower rewards.
- DD is not NE because D wants to change to C to get better reward.
- Although CD or DC is NE, but they are not stable because no one wants to always be C to receive a relatively lower reward.

- Mixed Strategy: Dare 1/3 of the time and Chicken out 2/3 of the time. the expected payoff is 14/3.



- Both agents know three cards [C,C], [CD], and [DC]
- A third person Chris choose a card randomly and inform the agent what to do base on the card. Should the agents listen to Chris?
 - If Chris told me to D, then other one must be C. I don't want to change because changing to C will get lower reward.
 - If I was told to C, then the other one has half the chance to be C and half the chance to be D. If I choose C, I will get expected reward as 4. If choose D, will get expected reward of 3.5 (see green words in the next slide). So staying with C is better.
- Yes.** And this is called a **correlated equilibrium(CE)**

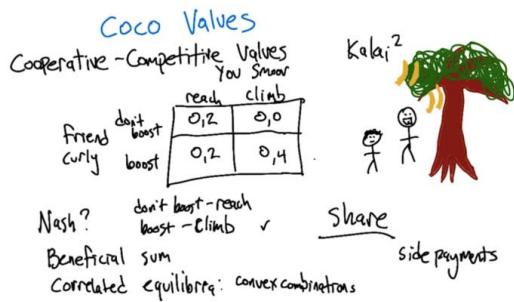


- Given [C,C],[CD], and [DC] are randomly chosen, the expected CE is $1/3 * 7 + 1/3 * 2 + 1/3 * 6 = 5$, which is better than 14/3.
 - Not both chicken out to get reward of 6?
Because [C,C] is not an equilibrium!
 - What are the real life analogy of "Chris"?
Traffic lights, stop signs.....

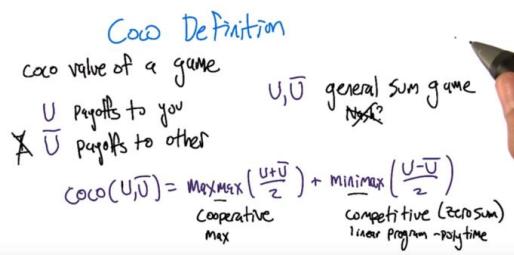
CORRELATED FACTS

- ① C.E. can be found in poly time
 - ② All mixed Nash are correlated,
so C.E. exist
 - ③ All convex combinations of mixed Nash
are correlated

Cooperative-Competitive Values



- There are two N.E in the CoCo game.
 - Curly are not motivated to do anything because he always gets 0 in the game. How to motivate him? Side payments? Smooth shares sth with him?



- To understand this, need to go back and revisit Maxmax, minimax and linear programming.
 - The basic idea is to rewrite CoCo into a cooperative part and a competitive part.

$$\begin{aligned}
 & \text{coco Example} \\
 U = & \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}, \bar{U} = \begin{bmatrix} 2 & 0 \\ 2 & 0 \end{bmatrix} \\
 \frac{U+\bar{U}}{2} = & \begin{bmatrix} 1 & 0 \\ 1 & 0 \end{bmatrix} \quad \frac{U-\bar{U}}{2} = \begin{bmatrix} 1 & 0 \\ -1 & -2 \end{bmatrix} \quad \frac{\bar{U}-U}{2} = \begin{bmatrix} 1 & 0 \\ 1 & 2 \end{bmatrix} \\
 & \text{maxmax} \qquad \qquad \text{minimax} \\
 & \text{U's perspective: coco value } \perp \qquad \bar{U}'s \text{ perspective:} \\
 & \text{side payments } P = \text{coco}(U, \bar{U}) - U(q^*, \bar{q}^*) = 1 \quad \text{below} \\
 & \bar{P} = \text{coco}(\bar{U}, U) - \bar{U}(q^*, \bar{q}^*) = -1
 \end{aligned}$$

- Given a U and U_{bar} matrix, compute the coco value and side payments.
 - Calculate maximax and minimax. There are minimax form U 's and U_{bar} 's perspective. But the two values cannot be values that's outside NE regions.
 - Coco value from U and U_{bar} 's perspective are different.
 - Side payments can be calculated based on U and U_{bar} 's max value and coco value.

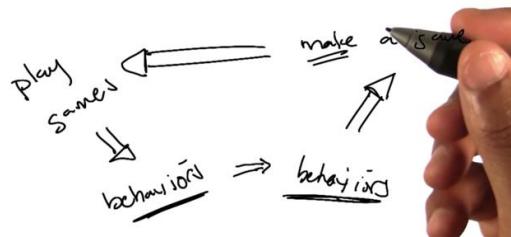
Coco Properties

- efficiently computable
 - Utility maximizing
 - decompose game into sum of two games
 - UNIQUE
 - Can be extended to stochastic games (coco-Q, not mechanism)
 - not necessarily equilibrium (best response)
 - needs to be binding always coco? Hiring police officers
 - Doesn't generalize n. 2. cococo

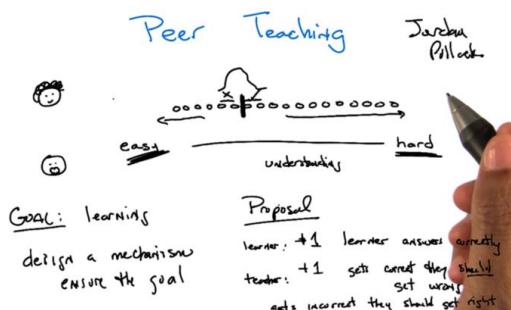
- CoCo-Q converges but it is not non-expansion.

Mechanism Design

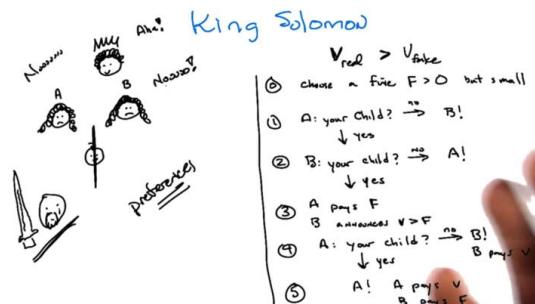
Mechanism Design



- Based on wanted behavior to build/make a game.
- E.G. Tax policy to encourage people buying a house



- Goal: learning.
- Learner get incentives by getting questions correct -> promote learning
- Teacher get incentives by setting the difficulty bar so that if the learning get an easy question wrong and get a hard question correct. -> push student to learn



- In the old King Solomon story, A and B will both say "no" and King will not get no real information.
- Given the mechanism on the left, fake mom will be honest.

- Assumption: the baby's value to real mother and fake mother will be $V_{real} > V_{fake}$
- Ask both the mother, is this your child, if all say yes, then A pays a fine F.
- Then B announce a V that is $V > F$
- Ask the question the same question to A again, if A say yes, A pays V, B pays F and A keeps the baby. If A says no, B keep the baby and B pays V.

Explanation: if B is fake, she will not announce a value that's larger than V_{fake} , because if A says no, he will get the baby with more than he is willing to pay. So A will get the baby. If B is real mother, she will announce a value V_{real} which is larger than V_{fake} . And A won't say yes because again she will end up with the bay but paying more than she is willing to. So, the real mom will get the baby.

What have we learned?

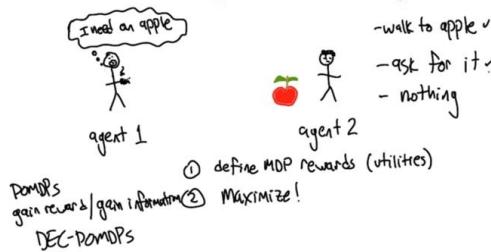
What Have We Learned ?

- Scooby Doo
- Solution concepts (Nash, Subgame Perfect)
- Correlated Equilibrium
- Mechanism design
- Peer teaching
- Baby division
- Subgame perfect
- Coarse values
- Nash
- Efficient
- Shared source of randomization
- Better outcome (?)
- Binding side payments
- Pareto optimal

Week 14: Coordinating Communicating and Coaching



Coordinating & Communicating



DEC-PoMDP

Decentralized Partially observable Markov Decision

I finite set of agents
 S states
 A_i agent i's actions
 T joint transition function
 R_i shared reward function (R_i if POSG)
 Z_i agent i's observations
 O observation function

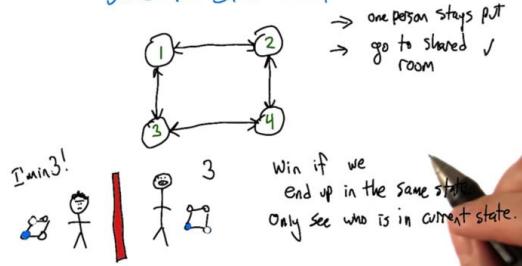
vector with one action per $|T|$
 $T(s, a, s')$

- Dec-POMDP has some perspectives of game theory and MDP
- Multiple agent working on getting a common reward. (if the rewards are separated for all the agents, then it's a POSG partially observable stochastic game)

DEC-PoMDPs

- Elements of game theory and POMDPs
- NEXP-complete (for finite horizon)
- Agents can benefit from communication
- optimal solution balance cost of communicating with cost of not communicating
- Algorithms, heuristics, applications known.

DEC-PoMDPs Example

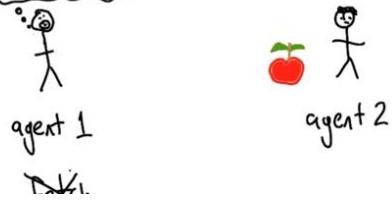


- Two agents, they know where they are but don't know the other's position. When the two are in the same room, they win.

- Strategy: go to a shared room. But my knowledge of my current position could be wrong (partially observable world).

Communicating & Coaching

I need Curly to learn to get an apple



- Agent 1 wants to set up some kind of reward function to move agent to do something (e.g. get the apple for me).

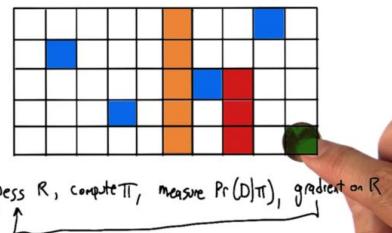
Inverse Reinforcement Learning

Inverse Reinforcement Learning

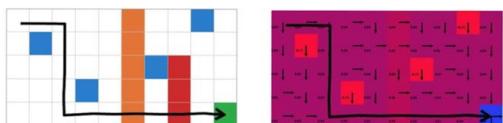
environment
 ↓
 rewards \Leftarrow IRL \Leftarrow behavior

- Inverse Reinforcement Learning: the agent experience the environment and a set a behavior and then generate a reward function based on the inputs.

Inverse Reinforcement Learning



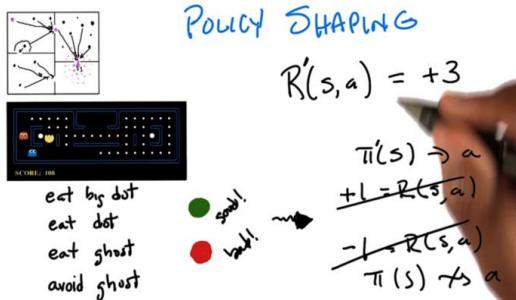
Output of MIRL



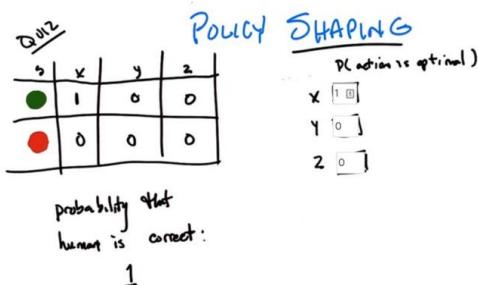
What Have We Learned?

Coordinating
Communication
Coaching

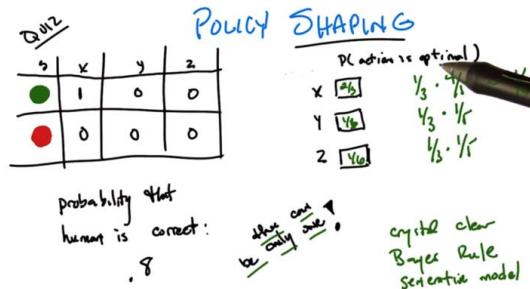
Policy Shaping



- If a human is giving feedback (commentary) about whether the agent's action is good or bad, s/he is doing policy shaping.
- Policy shaping could be realized by reward shaping which is replace reward of an action with a new reward?
- Agent need a mechanism to learn from the environment and the commentary to decide what policy to take (not just listening to the commentary, because the commentary might not be always right).



- If human is always correct, given the feedback, what's the probability that the action (x, y, or z) is optimal?
- Answers in the slides above.



- What if human is 0.8 probability of right?
- Counting method:
 - Saying x is optimal is like saying y and z is not optimal.
 - Since human is 0.8 correct, then x, y, z being optimal is 0.8, 0.2, 0.2.
 - Normalize the numbers above, will get 2/3, 1/6, 1/6.

Policy Shaping

$$p(a|d_a) = \frac{C^{\Delta_a}}{C^{\Delta_a} + (1-C)^{\Delta_a}}$$

General case

$\Delta_a = \# \text{ yes} - \# \text{ no}$

$C = \text{prob correct}$

only one! $p(a|d_a) \propto C^{\Delta_a} (1-C)^{\Delta_a}$

- Δ_a is coming from data of action a (d_a). C is the probability of correct of the people giving commentary.
- The formula above give the method of calculating probability of action a is optimal.
- Note: the final probability will need to be normalized against the probabilities of other actions.

Policy Shaping!

Multiple Sources!

Qwiz

	x	y	z
π_H	2/3	1/6	1/6
π_A	1/3	1/6	8/15
	o	o	o

- In the policy shaping case, information are coming from multiple sources.
- E.g. π_a and π_H are policy info from agent exploring the world and human giving feedback.
- Some algorithm decrease the importance of π_H as time goes. One need to know that π_a already incorporated the information of human uncertainty (C).
- The way to combine the two sources is to calculate the probability that the two policy will agree: $a_{opt} = \text{argmax}_a p(a|\pi_1) * p(a|\pi_2)$.
 - In the quiz $x_{opt} = 1/15$, $y_{opt} = 1/60$, $a_{opt} = 2/15$. So we should choose z as optimal.

Drama Management

Drama Management

- demonstrations
- rewards
- policies



Player ← Agent
Agent → Author

- The way a human can communicate to an agent
 - demonstration: show the agent what's the correct action (inverse RL)

- reward shaping: giving reward for agent's actions
- policy shaping: commentary on the agent's actions
- Author convey his intent to the agent so the agent can

Drama Management

What's a story?

- trajectory!
- though ... plot points

How can we influence a player?



- story can be defined as a trajectory through plot points

Trajectories as MDPs

states	: partial sequences	hyper exponential
actions	: story actions	
model	: player model	
rewards	: author evaluation	

- Above a some mapping of MDP elements to trajectory MDP elements
- Problems
 - large number of sequence of states (hyper exponential)
 - Since MDP will maximize rewards, treating story as an MDP will only make the author happy and force the player to experience the story.

Trajectories as TTDMDPs

stories
 trajectories : sequence so far linear T_t
 actions : ~~not~~ actions length T_t
 model : $p(t'|a_t, t)$ are player
 target distribution : $p(T)$
 policy : probability distribution over actions

- $p(t'|a_t, t)$ is the probability that the player at trajectory t and take action a then ended up in trajectory t' . $P(T)$ is a target distribution.
- The action is not player's action but the story action
- The optimal policy is the policy that will lead to the targeted trajectory distribution $P(T)$
- The calculation time is linear and dependent on the length of the story.

What have we learned?

WHAT HAVE WE LEARNED?

- ④ Combinations, communicating, coaching
- DEC POMDP
 - shared $R_t(a)$
 - distinct brains
- ⑤ IRL
 - behavior → reward
 - reward shaping
- demonstrations
- policy shaping
 - TTD-MDPs
 - constraints