**RL 1 - Markov Decision Process**
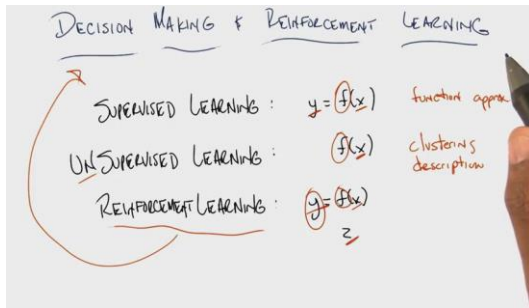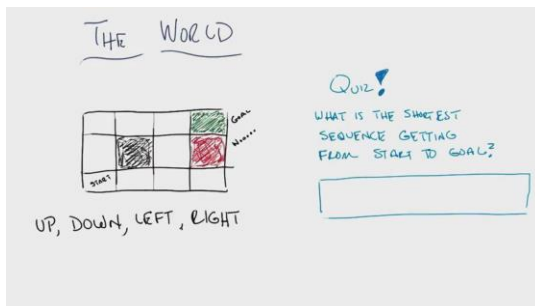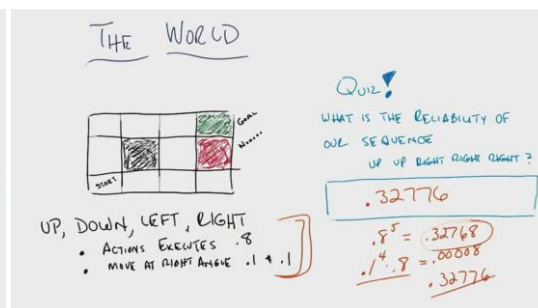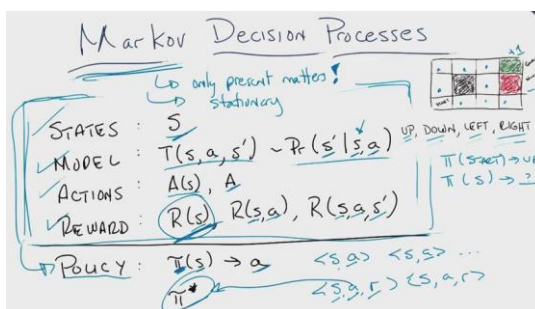


- For reinforced learning, given x and z, and try to find function f that generates y



- U, U, R, R, R or R, R, U, U, R
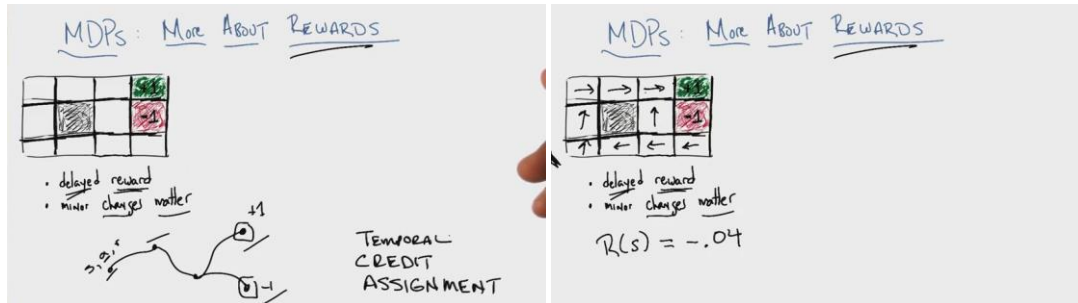- Usually multiple solutions available



- Stochasticity of movements
- Need to consider the path to the right
- Need to finish in green spot after 5 steps, no more steps available
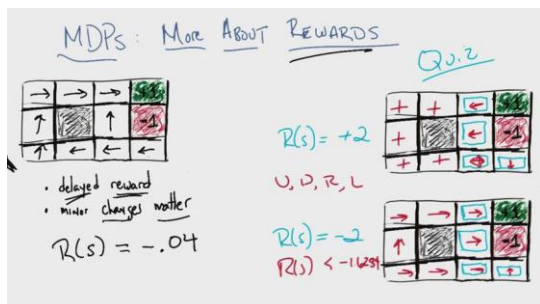


- add historical states to make the process Markov
- RL: given <state,action,reward> pairs as training samples, and try to find the optimal policy pi or

optimal solution

- SL: given <state,action> pairs as training samples, and try to find the optimal polity pi
- Assumptions: you know which state S, and the reward for taking actions
- Policy does not tell you the sequency of actions (i.e. plan) to reach success, but what actions to take at particular state
- The problem: given the MDP, how to find the optimal policy pi?



- MDPs: more about rewards
- Delayed rewards
- Temporal credit assignment, where credit = 1 or -1
- reward R(s) = -0.04 for each state except the terminating positions
- The arrow the global optimal solution (or plan)
- The reward function may change the global solution greatly: totally different path given different reward functions



- Big negative rewards to end game as soon as possible
- Big positive rewards to continue game as long as possible
- Medium negative rewards: in the game to get the +1 reward
- Big positive rewards: just in the game
- Big negative rewards: get out of game ASAP
- Reward function is domain knowledge: if you want to use MDP to describe a world, you need to think carefully how you set up reward functions

- Infinite horizon: you can live for ever to try, till get the terminating states => need to consider reward function, as well as how much time left (i.e. horizon)

- If the horizon is not infinite, for example, game ended within 3 steps when in state (3,1), then probably going UP is the right choice, instead of going left

- In the infinite horizon, the process is stationary, where the action only depends on the state. Or, as long as I am at this particular state, I will always take the same action. Policy pi(S) -> a

- In the finite horizon, the process is non-stationary, where the action may be different depending on how much time left at the same state. Policy pi(S,t) -> a

- Utility of sequence (vs utility of a single step): $U(s_0, s_1, ...,)$

- Lemma (stationarity of preference): if $U(S0,S1,S2,...) > U(S0,S1',S2',...)$ (utility of sequence starting at the same state s0), then

- $U(S1,S2,...) > U(S1',S2',...)$

- Or, if I prefer a sequence over another sequence today, then I prefer the same sequence tomorrow

- Note: the utility of a sequence = sum of rewards along the path, then the Lemma is true. On the other hand, if Lemma is true for any sequence, then the utility of sequence = sum of rewards along the path.



- Discounted rewards by adding the factor gamma^t
- Infinite horizon but finite total reward

POLICIES

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

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

immediate    long term    delayed reward

POLICIES

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

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

immediate    long term    delayed reward

$$\pi^*(s) = \arg\max_a \sum_{s'} T(s, a, s') U(s')$$

$$[U(s) \equiv U^{\pi^*}(s)] \quad U(s) = R(s) + \gamma \max_a \sum_{s'} T(s, a, s') U(s')$$

Bellman Equation

- The optimal policy pi, will optimize the total expected discount rewards
- The utility given policy pi is the expected total discount rewards given policy pi and starting state S
- T(S,a,S') is the transition probability that we start at state S and end up at state S' under activation a
- U(S) is the total discount rewards given optimal policy
- Bellman equation is the key to solve MDP problem



POLICIES : FINDING POLICIES                    not-linear

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

n equations    in    n unknowns

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

$$\hat{U}_{t+1}(s) = R(s) + \gamma \max_a \sum_{s'} T(s, a, s') \hat{U}_t(s')$$

POLICIES : FINDING POLICIES                    not-linear

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

n equations    in    n unknowns

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

Value Iteration

$$\hat{U}_{t+1}(s) = R(s) + \gamma \max_a \sum_{s'} T(s, a, s') \hat{U}_t(s')$$

truth

- For Bellman equation, there are N equations with N unknowns, but max function make it not-linear => algorithm is used to solve these equations (called value iterations)
- U_t(S) is arbitrary values, while R(S), the instant reward function, is the truth
- Adding more truth through loop will converge to the solution
- Value iterations: does not give you the right answer, but closer and closer answers
- Once the utility function is solved, then it's easy to solve the optimal policy pi



POLICIES : FINDING POLICIES

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

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

QUIZ!

$$-.04 + \frac{1}{2}[0 + 0 + .8] \quad U_1(x) = \boxed{.36}$$

$$-.04 + \frac{1}{2}[.036 + -.004 + .5] \quad U_2(x) = \boxed{.376}$$

$$\gamma = \frac{1}{2}, \quad R(s) = -.04, \quad U_0(s) = 0 \quad \pi(s) \Rightarrow a_{utilities}$$
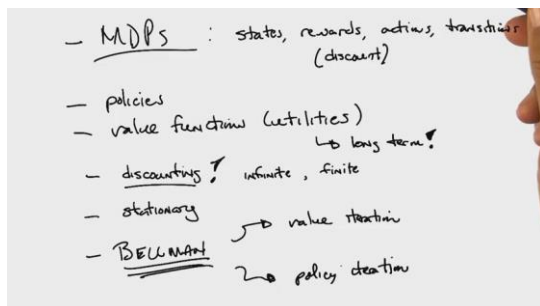
- For state X, better to go to green spot with probability 0.8, while going to other two directions with probability 0.1
- For t=2, we need to calculate the U(Y) at t=1 where Y is the grid below X: the optimal policy is to go to the wall, and back to the original state. However, at some point, when the utility of state X is big enough, the optimal policy would be go to state X, instead of wall
- The value iteration works if eventually the value pops out its neighbors => +1 and -1 will pop out to

states over time

- Given utility, we can find policies; but given policies, it's hard to find utility
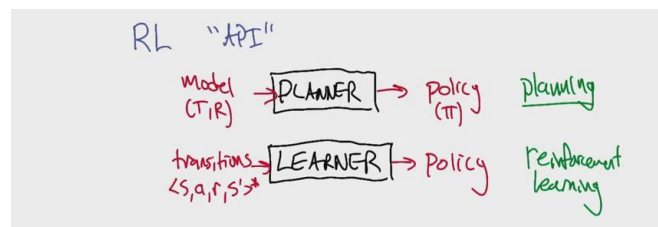


- Given policy, we can remove the max function, and turn the Bellman equation into linear equation, and solve it
- Called policy iteration
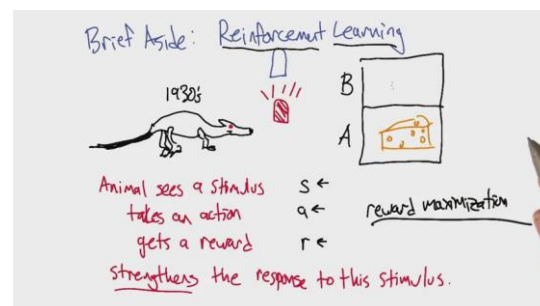- Guaranteed to converge: finite number of policies, and you are always getting better



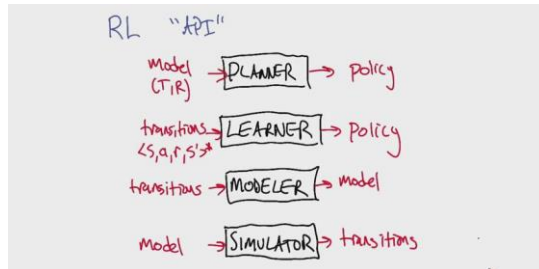- In RL, you don't necessarily know the rewards and transitions
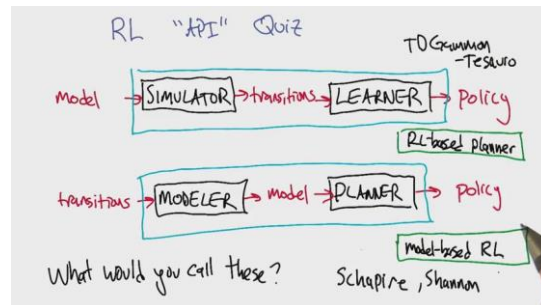
**RL 2 - Reinforcement Learning**



- T: transition matrix; R: reward functions
- \pi: S-> actions; Transitions: samples

- Stimulus: like seeing the red light  <- state; Action: like picking box A    <- action
- Reward: like cheese    <- reward; Strengthen: means next time sees red light, more likely to go to box A => reinforcing the connection
- How to choose actions to maximize reward as a function of state   => reward maximization



- Modeling: like machine learning problem
- Simulator: simulate samples, time consuming
- One way to do learning: to do modeling and simulating, so that I know a model and what the reinforcement function was, and just do planning





- 3 outputs of RL algorithms: policy search algorithm; value-function based algorithm; model-based RL learner
- From utility to policy: use argmax operation
- From model to utility: value iteration to solve Bellman equation
- Model based RL learning: direct learning, not-direct usage
- We will focus on the second: value-function based learning: medium direct learning and usage

A new kind of value function

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

$$\pi(S) = \text{argmax}_a \sum_{s'} T(S,a,S') U(S')$$

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

Value for arriving in S, leaving via a, proceeding optimally thereafter.

A new kind of value function

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

$$\pi(S) = \text{argmax}_a \sum_{s'} T(S,a,S') U(S')$$

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

$$U(S) = \boxed{\max_a Q(S,a)}$$

$$\pi(S) = \boxed{\text{argmax}_a Q(S,a)}$$

Use Q to define U & π.

- Q(S,a) is the reward taking a specific action => we can take Q as U, if we always take the best action. The best action is the one maximize the utility (value you get from that point on)



Q-learning : The Quiz

o figuring out the best line to wait in
   queue-learning
o discovering when to come in
   for your line
   cue-learning
o practicing the best bank shot
   cue-learning
o evaluating the Bellman equations from data
   Q-learning

Estimating Q from transitions

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

transition
<S,a,r,S'>:

$$\hat{Q}(S,a) \xleftarrow{\alpha} r + \gamma \max_{q'} \hat{Q}(S',q')$$

Q-learning update equation

utility of next state
utility of state

V & X
V ← (1-α)V + αX

- We don't have access to R(S) and T(S,a,S'). if known, we can use value iteration and policy iteration to solve Bellman equation for Q(S,a)
- <S, a, r, S'>: we are in state S, take the action a with reward r. The new state is S'
- V<-alpha X: current state is V, and move alpha toward state X
- V<- (1-alpha)V + alpha X



Learning incrementally

$$V \xleftarrow{\alpha_t} X_t$$

$$X_t \sim X$$

What does V Converge to?
o E[X]
o it doesn't
o var(X)
o ∞

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

$$\alpha_t = \frac{1}{t} \quad \sum_{t=1}^{t} \frac{1}{t} \approx \ln t$$

$$\sum \frac{1}{t^2} = \frac{\pi^2}{6}$$

Learning incrementally

$$V \xleftarrow{\alpha_t} X_t$$

$$X_t \sim X \quad \text{iid}$$

What does V Converge to?
✓ E[X]
o it doesn't
o var(X)
o ∞

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

$$\alpha_t = \frac{1}{t} \quad \sum_{t=1}^{t} \frac{1}{t} \approx \ln t$$

$$\sum \frac{1}{t^2} = \frac{\pi^2}{6}$$

Average!!

- The X_t samples from the same distribution X, IID
- The learning rate alpha_t satisfies the following conditions: sum is infinite but sum of square is finite

Estimating Q from transitions

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

$\langle S,a,r,S' \rangle:$

$$\boxed{\hat{Q}(S,A) \xleftarrow{\alpha} r + \gamma \max_{a'} \hat{Q}(S',a')} \quad \text{changing over time}$$

$$= \mathbb{E}[r + \gamma \max_{a'} \hat{Q}(S',a')]$$

$$= R(S) + \gamma \mathbb{E}_{S'}[\max_{a'} \hat{Q}(S',a')]$$

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



Q-learning convergence

$\hat{Q}$ starts anywhere

$\langle S,a,r,S' \rangle$ $\quad \hat{Q}(S,A) \xleftarrow{\alpha} r + \gamma \max_{a'} \hat{Q}(S',a')$

then $\hat{Q}(S,A) \to Q(S,A)$. Solution to Bellman equation!

if $S,A$ visited infinitely often ←

$\sum_t \alpha_t = \infty$ , $\sum_t \alpha_t^2 < \infty$

$S' \sim T(S,A,S')$ , $r \sim R(S)$

- Q learning:
- Under alpha learning, the right side will go to the expectation over time
- The green formula derivation is not quite right, as Q hat is changing over time
- In order to be convergent, need to visit states for infinite times



Choosing Actions

Q-learning is a family of algorithms.
→ how initialize $\hat{Q}$?
→ how decay $\alpha_t$?
→ how choose actions? — always choose $a_0$ (won't learn)
— choose randomly (won't use it)
— use Q (won't learn it) (won't learn)

$\forall S \; \hat{Q}(S,a_0) =$ awesome (chicken dollars)  greedy
$\forall S, a \neq a_0 \; \hat{Q}(S,a) =$ terrible   "local min"

- Q learning is a family of learning algorithms. It differs by 3 dimensions



Choosing Actions

Q-learning is a family of algorithms.
→ how initialize $\hat{Q}$?
→ how decay $\alpha_t$?
→ how choose actions? — always choose $a_0$ (won't learn)
— "simulated annealing"-like approach >— choose randomly (won't use it)
take a random action sometimes >— use Q (won't learn it) (won't learn)
$\hat{\pi}(S) = \arg\max_a \hat{Q}(S,a)$ w.p $1-\varepsilon$ — random  greedy
random action  otherwise$(\varepsilon)$  restarts! "local min"
(slow!)



$\varepsilon$-Greedy Exploration

If GLIE (decayed $\varepsilon$)
"greedy limit + infinite exploration"

$\hat{Q} \to Q$ and $\hat{\pi} \to \pi^*$ !!
learn  use

Exploration - Exploitation dilemma

one of you!

Fundamental tradeoff in RL

RL
model learning ↓ planning
ML  KWS

model-based
knows what you know

- Simulated annealing: will random act with probability epsilon: will explore to all states



What have we learned?
— learn to solve an MDP (T,R), interact $\langle S,a,r,S' \rangle$
— Q-learning: converge, family  initialize $\hat{Q}$
set $\hat{Q} =$ awesome
— Exploration-exploitation: learn & use!

— approaches to RL  optimism in the face of uncertainty
— connection to planning
$A^*$  → function approximation generalizing

**RL 3 - Game Theory**



- Game theory can be thought about in different perspectives, like listed

- Number is what assigned to a

- a got two actions: L and R, while b has 3 actions: L, M, R

- zero-sum: the sum of rewards is constant (not necessarily zero)

- no statistical transitions

- **strategies** (in game theory): mapping from state to actions (similar ideas as policy in MDP setting), always refer to one party (like a)



- 4 choices for a, even though if I go R at stage 1, no need to make choice at state 4. But independently, there are 4 choices

- Strategies: what you might do in all states you might end up with



- Form a matrix, whose value for A is the combined actions of A and B

- Once we write down the matrix, then nothing else matters. Then the graph can be removed. Everything about the game is captured in the matrix

- This is called the matrix form of the game. The L, R, M combinations do not matter as well. What matters is by following strategies, you got values XXX for A, and how you get there, does not matter

- RL is to optimize long-term reward -> pick a policy to have a best long term reward possible
- Here a has 4 strategies, while b has 3 strategies. Based on this matrix, what will a, and b take?
- The same minimax answer if we choose row first or column first under 2 player, zero-sum, deterministic perfect information game
- If you have the matrix, you can construct the tree based on it (may have multiple tress corresponding to the same matrix)



- Minimax theory
- Assuming the other party take the counter strategy
- The value of this game for a is 4 if both players are rationale
- a first or b first, the results are the same
- Assumption: everyone is trying to maximize reward



- Square box: stochastic rewards; Circle box: deterministic rewards
- The randomness does not necessarily happen to leaf only under general non-deterministic game
- We can reconstruct a tree based on the matrix, but infinitely many trees corresponds to the matrix
- The minimax theorem still holds tree for non-deterministic game (Von Neumann theorem)
- The only thing matters is the matrix, and use minimax or maxmini to figure out the solution



- Red is bad for A, while black is good for A

- Color of A is random, with 50% to be red and 50% to be black

- B doesn't know the result of A

- Hidden information, as B does not know the result of A

- A will not resign for black card

- B has not idea which state to be in (figure above right), hence hidden information

- A hold or resign based on red card (sure strategy for black card)

- B will resign or see when A choose to hold (2 strategies)



- A resigner and B resigner, means when A got red, resign ->-20. But when A got black, always wins -> +10 => 0.5*(-20 + 10) = -5

- Hidden information, means the order of making choice make differences

- If A chooses first, -5; if B chooses first, +5

- If one is consistent over his actions, then the other party can manipulate it and earn profits <- why order of choosing matters

- Pure strategy (or consistent strategy) -> mixed strategy (distribution over strategies. I.e. 30% holder, and 70% resigner)



- Quiz 1: B always resigner, and what's A's expected profit given probability of being A holder

- Two lines corresponding to the two strategies by B
- If A choose a mixed strategy with probability 0.4 choosing holder, and 0.6 choosing resigner, then the expected payoff is +1 to A, irrespective of what B chooses
- If B chooses to use mixed strategy as well, the final payoff will be the same as +1 to A
- If B chooses to use mixed strategy, payoff would be somewhere between these two lines
- The highlighted triangle (above right) is the payoff of A choosing different probabilities, and we choose to take the maximum
- If the triangle is decreasing (see the cross above), we will choose the left end point, instead of the cross point as the solution
- If B choose to use mixed strategy, and then consider A, the results would be the same as if we let A choose mixed strategy, and then consider B



- For non-zero sum game, need to put down the payoff for both players, instead of one number
- Based on the matrix, the best solution is both corporate and sit in jail for 1 month each
- Not joint choice, but individual choose
- But actually ended with both defect
- If A know B is going to corporate, then A should choose defect to save 1 month;
- If B know A is going to defect, B should defect to save 3 months
- Eventually end up with both defect
- For each column, A is always better to defect than corporate (0 vs -1 and -6 vs -9), so A will never corporate



- In (pure) Nash equilibrium, no player would like to change their strategies for better payoff (assuming all the other keep constant)

- In the first case, the second row strictly dominate the first row
- In the second case, both have the largest payoff



- Listed are 3 fundamental theorems for Nash Equilibrium
- 1st: in iterative fashion: it may not dominate in one iteration. Could do row&col iterations
- 3rd: for finite game, there exist at least 1 Nash Equilibrium (could be mixed)



- Get communication instead of wall may not solve the prisoner dilemma, because the person goes first, the person lose
- This happen only when the game plays once. If game happens more times => keep corporating together for ever; corporating will earn reward
- Suppose we play the game 2 times. For A, either we can set up like 4 strategies: either corporate or defect at each round; or we can let A's choice depends on B's choice in the first round, ending with 8 strategies
- For example, we play the game for 20 times. If we go forward, it might make sense to corporate. But going backward, at 20th game, whatever happen is the sunk cost. Final one is determined (-6, -6). At 19th game, 20th game result is there, so 19th game will both defect as well => even if we can play for many times, we will still end up at the same solution
- **Theorem**: for N repeated games, the solution is N repeated Nash Equilibriums. If a game has multiple N.E. point, the results could oscillating => once you are in one N.E. point, we are not going to leave it

- **Mechanism Design**: to solve the prisoner dilemma, change the game! (change utility) 1) If A cares 50% of B staying the jail, then the payments shift; 2) the people defect, got punished. 3) hire policemen or hitters

**RL 4 - Game Theory Continued**





- Conclusion: if the number of rounds is unknown, the results would be different

- Gamma here is the same as the discounting factor gamma





- In uncertain number of rounds, we cannot write down strategies like C-D-C-D-D, or in the form of trees. We need some representations for unbound round of games

- Tit for Tat strategy





- Listed are the expected payoff when we are dealing with Tit for Tat: always defect or always

corporate

- For gamma < 1/6, we should always defect; while gamma > 1/6, we should corporate

- The matrix form is good if we only play once

- Use MDP to solve the strategy: mapping from state to action => only 3 results (policies): always corporate (good for high discount factor), always defect (good for low discount factor), or D-C-D-C...

- If this is MDP, then there is no history. When you are in C, there are only two choices: C or D; the remaining will require you remember what you have done steps ago

- MDP always has a Markov deterministic policy



- Note: we are not playing with Tit for Tat

- There are two Nash Equilibrium points: one of them is corporative



- Folk theorem: everybody can prove it, but not sure who proved it first

- For quiz: just to execute some joint strategies, and see the payoff. Not necessarily to be in Nash equilibrium



- In the game, B and S tends to get the other lower score

- Smooth may choose mixed decision

- Minimax strategy = 1, while mixed strategy = 2/3



- Feasible region: possible payoff by using mixed strategies

- Acceptable region: both smooth and curly can get better than worst results



- Another perspective to prove Folk theorem is to use Grim Trigger

- Corporate (mutual benefit) until the cpty cross the line, and perform revenge for ever

- If A know B is taking the strategy above, then A will probably hold mutual benefit => no one has incentive to cross the line

- Problem: in some case, the threats is implausible (too crazy, nobody will do it)
- Plausible treats will cause subgame perfect: each player always takes best response independent of history
- Could take a better action based on its strategy, given some history
- TFT vs TFT is not subgame perfect: suppose we start with (D,D), then it will go with (D,D) for ever. We can step in and change it to (C,C) to have a better payoff



- We can check four combination and see whether they are subgame perfect



- MDP & RL is like stochastic game & multi-agent RL
- Dash line is the semi-wall, meaning if A would like to go through it, 50% will go through and 50% will stay put

- Goal: to get dollar
- They cannot occupy the same place; once one player reach the dollar sign, the game is over
- If they both go to the same position, they will throw a coin, and decide who will occupy the position
- When B change the strategy from going north to going west, then A's chance stay the same, while B's chance decreases. This is Nash equilibrium. The same for A going east, while B going north
- Both going to center is not Nash equilibrium, as you can go north



- Stochastic game setup



- Markov if the results do not depend on the other player
- Only one state in 3



- The maximum is taken under the assumption that everyone is trying to benefit me => works for team game, or somebody is willing to sacrifices themselves => therefore change max to minimax



- Minimax can be solved using linear programming
- Getting the second agent does not change much the problem
- One agenda: solved in polynomial time; while in stochastic game, it's not known to solve using

polynomial time



- Generate the framework to general stochastic game (not zero sum), then minimax operator will be replaced by Nash, as minimax does not suite the general case

- Value iteration does not work: if you loop through, weird thing will happen; does not necessarily converge and solve the equation

- Different Nash equilibrium may have different values => no unique solution

- Nash is defined using joint behavior, cannot be computed independently

- Two half Nash equilibrium cannot be Nash equilibrium => incompatible

- PPAD is known to be similar to NP, very hard to solve

- Based on Q function, cannot find what to behave => insufficient information



- Cheap talk: give the opportunity to let player to corporate a little bit

- Correlated equilibrium: version of Nash equilibria, more efficient to compute

- Give some of the reward back to you => to encourage corporate



- PD: Prisoner Dilemma