

Lecture5 Markov Decision Processes

Problem Many applications are non-deterministic! Go to next states with a probability

- Hidden Factors
- No need to consider all of them
 - Too complicated to consider (Too Many)
 - Not everything can be clearly understood

- Applications
- Robotics Decide where to move, but actuators can fail etc
 - Resource Allocation Decide what to produce,don't know customer demand
 - Agriculture Decide what to plant, don't know weather

Markov Given the present state, the future and the past are independent

- Markov Process
- Markov Reward Process (No Control)

$$P(s_{t+1} | \overbrace{s_t, s_{t-1}, \dots, s_0}^{\text{future current past}})$$
$$= P(s_{t+1} | s_t)$$

where t is the time
 - Type
 - Markov Decision Process (Control)

$$P(s_{t+1} | \overbrace{a_t, s_t, s_{t-1}, \dots, s_0}^{\text{future current past}})$$
$$= P(s_{t+1} | a_t, s_t)$$

where t is the time

- Markov Decision Process
- S State & S' New State After Action
 - Action(s)
 - $T(s, a, s')$ Transition Function $P(s' | s, a)$
 - $R(s, a, s')$ Reward Function $E(e | s, a, s')$

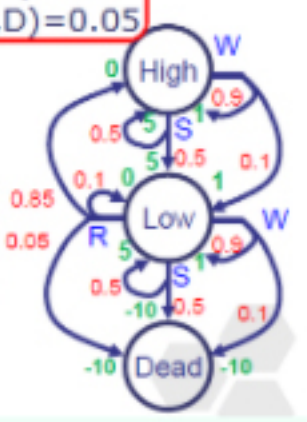
Recycling Robot

- $S = \{\text{High, Low, Dead}\} = \{H, L, D\}$
- Action(s) Action(H) = {Search, Wait} = {S, W}
Action(L) = {Search, Wait, Recharge} = {S, W, R}
Action(D) = {}
- $T(s, a, s')$

$T(L, S, H)=0.0$	$T(L, W, H)=0.0$	$T(L, R, H)=0.85$
$T(L, S, L)=0.5$	$T(L, W, L)=0.9$	$T(L, R, L)=0.1$
$T(L, S, D)=0.5$	$T(L, W, D)=0.1$	$T(L, R, D)=0.05$
$T(H, S, H)=0.5$	$T(H, W, H)=0.9$	
$T(H, S, L)=0.5$	$T(H, W, L)=0.1$	
$T(H, S, D)=0.0$	$T(H, W, D)=0.0$	
- $R(s, a, s')$

$R(s, R, s')=0$, except $R(s, R, D)=-10$
$R(s, S, s')=5$, except $R(s, S, D)=-10$
$R(s, W, s')=1$, except $R(s, W, D)=-10$

s and s' is any state, a is any action



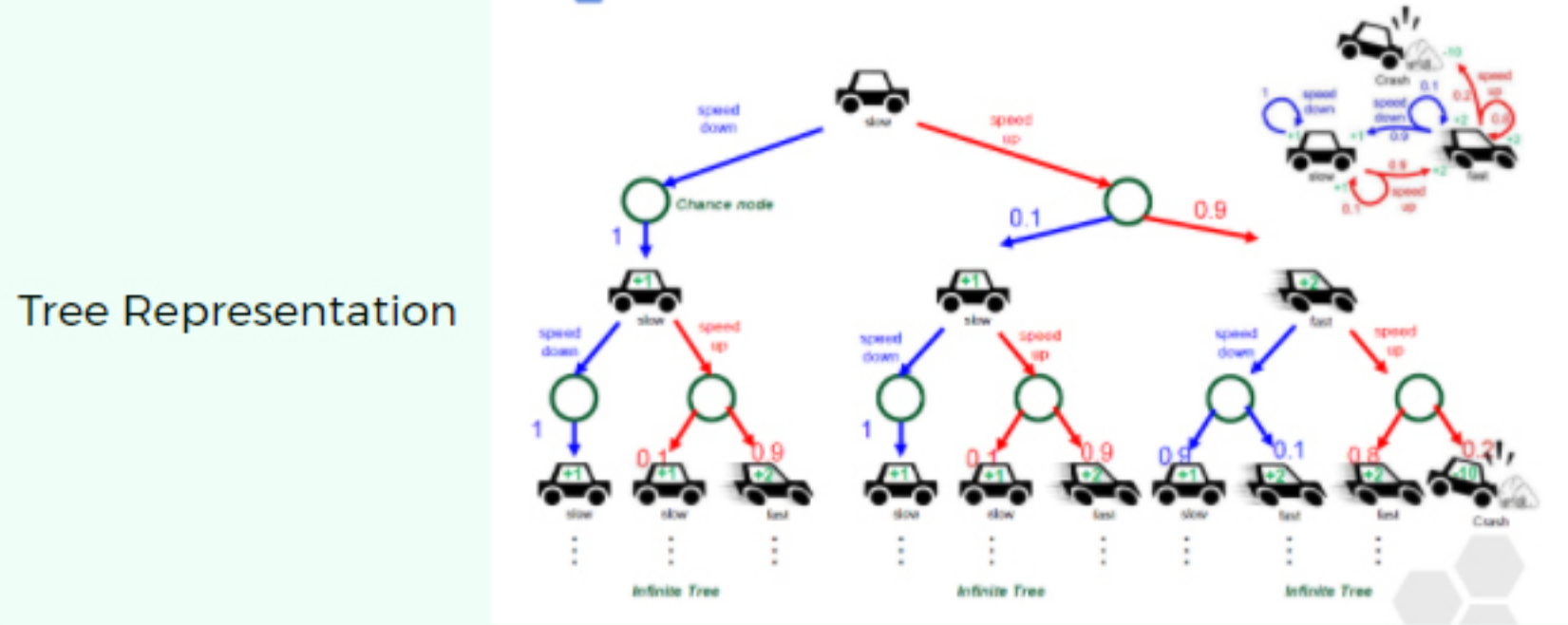
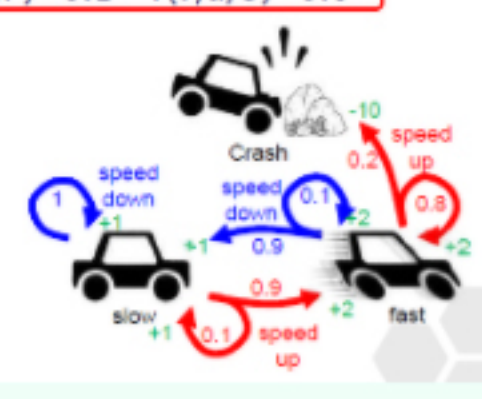
Car Driving

- $S = \{\text{Slow, Fast, Crash}\} = \{S, F, C\}$
- Action(s) = {Speed Up, Speed Down} = {u,d} s is any state
- $T(s, a, s')$

$T(S, u, S)=0.1$	$T(S, u, F)=0.9$	$T(S, u, C)=0.0$
$T(F, u, S)=0.0$	$T(F, u, F)=0.8$	$T(F, u, C)=0.2$
$T(S, d, S)=1.0$	$T(S, d, F)=0.0$	$T(S, d, C)=0.0$
$T(F, d, S)=0.9$	$T(F, d, F)=0.1$	$T(F, d, C)=0.0$
- $R(s, a, s')$

$R(s, a, S)=1$
$R(s, a, F)=2$
$R(s, a, C)=-10$

s is any state, a is any action



Problem: Tree can be extended forever and Actions may have infinite rewards

- Reward Estimation
- Time-Limited Values T Limit the maximum steps in a episode and end
 - Rewards decay exponentially

$$V(s_t) = \underbrace{R_{t+1}}_{\text{now}} + \underbrace{\gamma R_{t+2}}_{\text{Next}} + \underbrace{\gamma^2 R_{t+3}}_{\text{Next Next}} + \underbrace{\gamma^3 R_{t+4}}_{\text{Next Next Next}} + \dots$$

- Discount Factor
- $\gamma^{\infty} = 0$ means no influence to the result
 - Bellman Equation

$$v_{\pi}(s) = \sum_a \pi(s, a) \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma v_{\pi}(s')]$$

- Markov Decision Process Policy
- Prediction Given policy π , estimate state value $v_{\pi}(s)$
 - Control Estimate the optimal state value $v_{\pi^*}(s)$ pi * implicitly obtained

- State Value Estimation
- Value Iteration 根据policy迭代计算V(s)从0-n, 直到收敛converge
 - Policy Iteration
 - Calculate utilities for the fixed policy until converge

$$V_{k+1}^{\pi_i}(s) = \sum_{s'} T(s, \pi_i(s), s') [R(s, \pi_i(s), s') + \gamma V_k^{\pi_i}(s')]$$
 - Update the policy

$$\pi_{i+1}(s) = \arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^{\pi_i}(s')]$$
 - Policy Improvement
 - Loop until no update on policy
 - Comparison
 - Value Iteration
 - Update the value (and the policy implicitly)
 - Don't track the policy explicitly
 - Policy Iteration
 - Update utility in a fixed policy
 - After one policy is evaluated, a new policy is chosen
 - Both of them can provide the optimal solution for MDPs